



Introduction to Recurrent Neural Network

**Prof. Chia-Yu Lin
Yuan Ze University
2021 Spring**



Outline

- Industry 4.0 problem
- Why RNN?
- RNN Model
 - Recurrent Neural Network (RNN)
 - Variants of RNN
 - Long Short-term Memory (LSTM)
- Learning in RNN
- More applications of RNN
- RNN for anomaly detection
- Summary



Why RNN



Can machine understand the meaning of words?

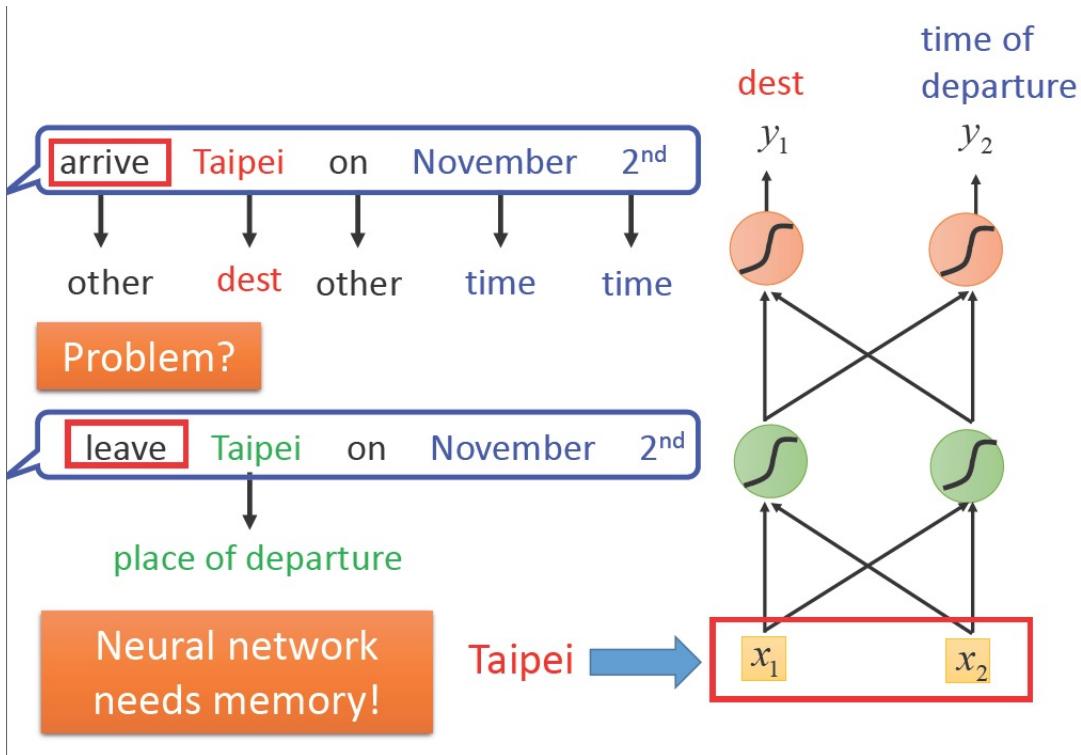
“Dog” is close to “Cat”

“Love” is close to “Hate”

“Korea” is close to “America”

Why RNN?

- Position of words is important
- Slightly change a word may change meaning of whole sentence.
- Lack of sequence concept in ordinary NN.



Memory is important

	x^1	x^2	x^3		
Input: 2 dimensions	4 7	4 7	1 1		
				1 1	
				1 4 4	
				+ 1 7 7	
					<hr/>
					3 2 1

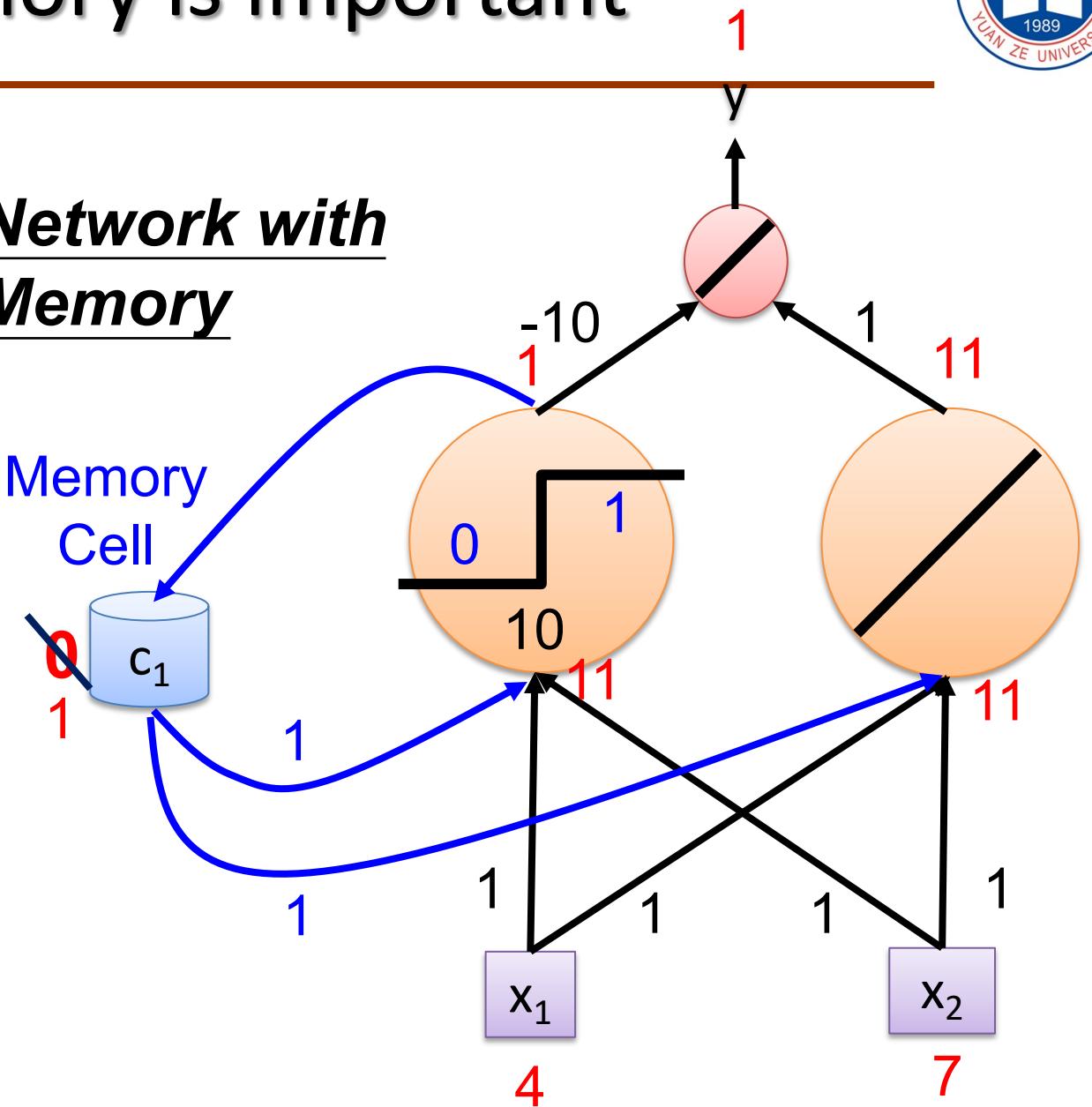
	\hat{y}^1	\hat{y}^2	\hat{y}^3		
Output: 1 dimension	1	2	3		

Network needs
memory to achieve this

Memory is important

Network with Memory

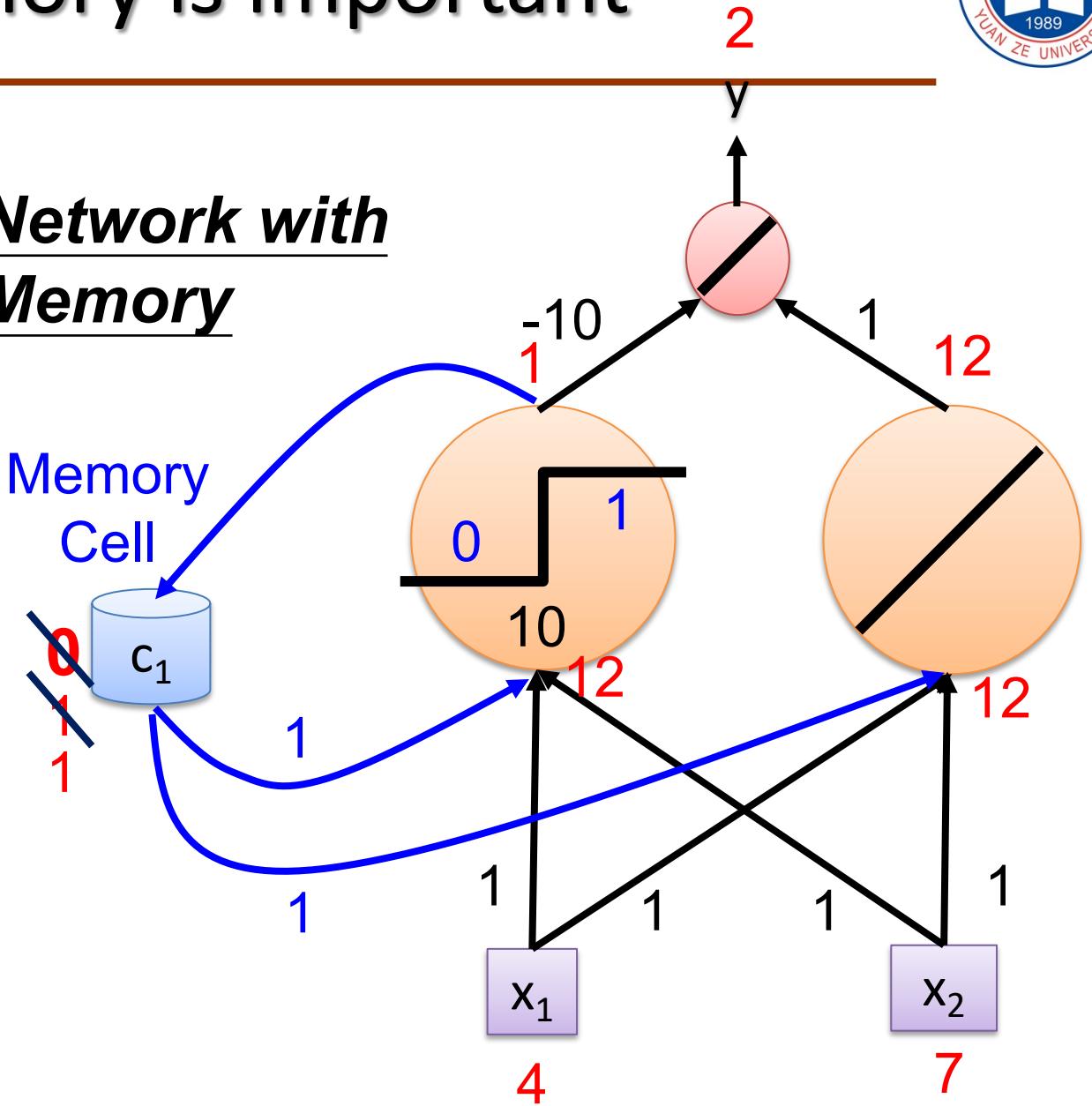
x^1	x^2	x^3
4 7	4 7	1 1
\hat{y}^1	\hat{y}^2	\hat{y}^3
1	2	3



Memory is important

x^1	x^2	x^3
4 7	4 7	1 1
\hat{y}^1	\hat{y}^2	\hat{y}^3
1	2	3

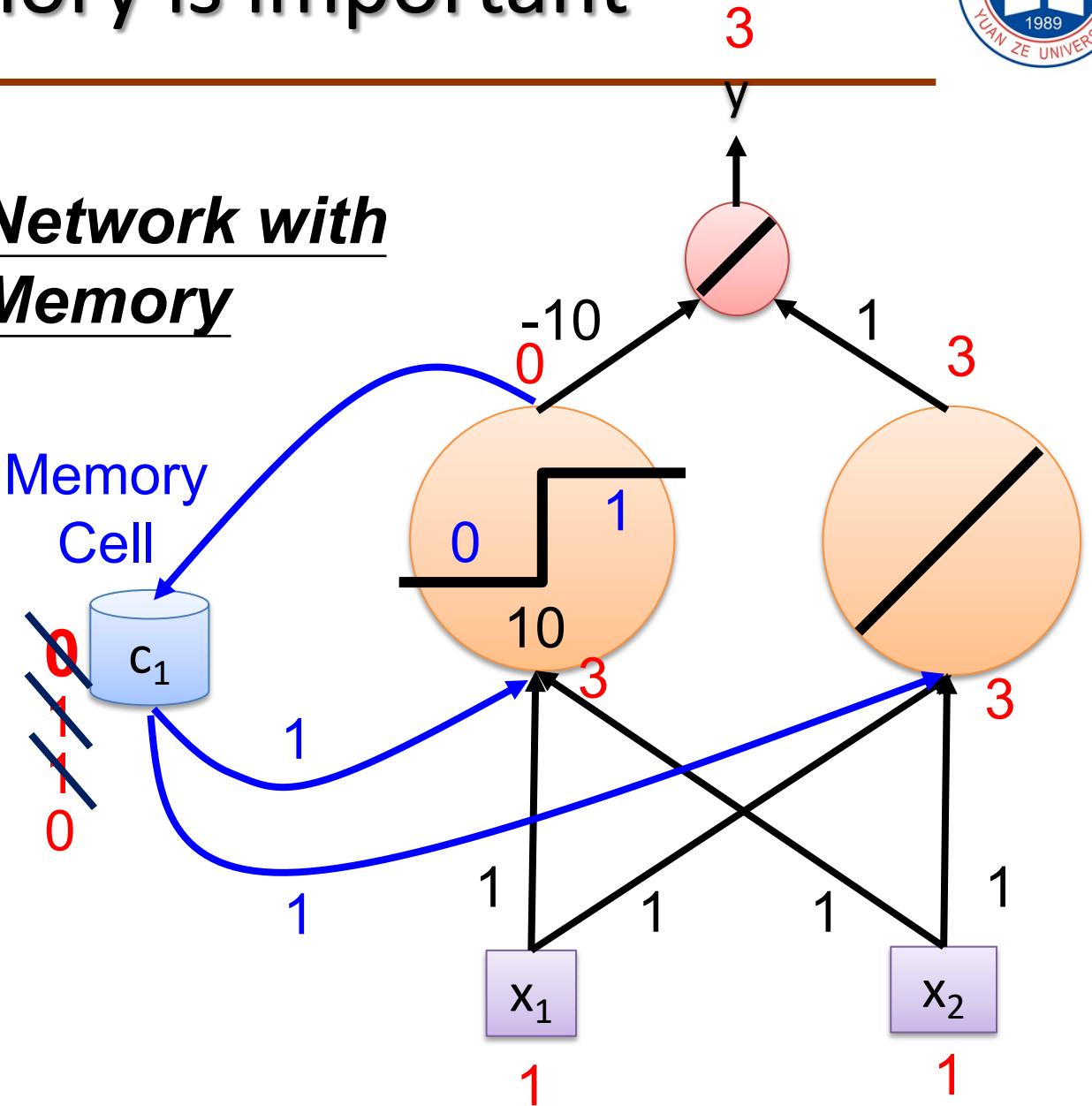
Network with Memory



Memory is important

x^1	x^2	x^3
4 7	4 7	1 1
\hat{y}^1	\hat{y}^2	\hat{y}^3
1	2	3

Network with Memory





RNN Model



Outline

Recurrent Neural Network (RNN)



Variants of RNN



Long Short-term Memory (LSTM)



Outline

Recurrent Neural Network (RNN)



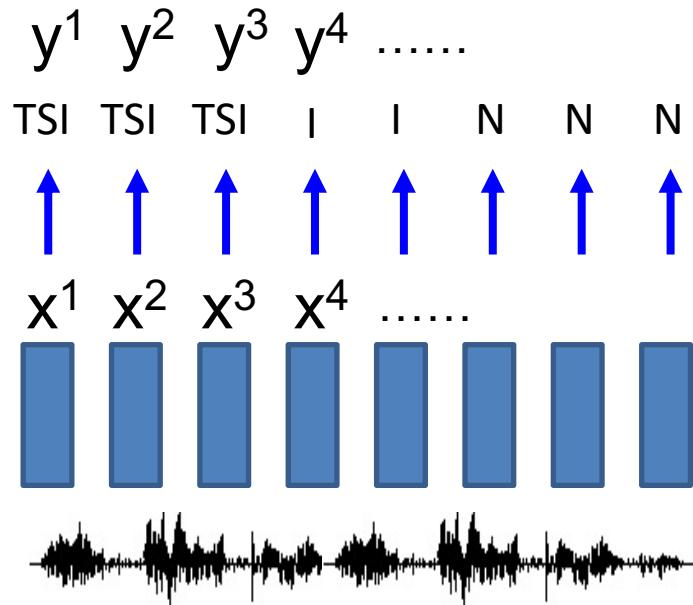
Variants of RNN



Long Short-term Memory (LSTM)

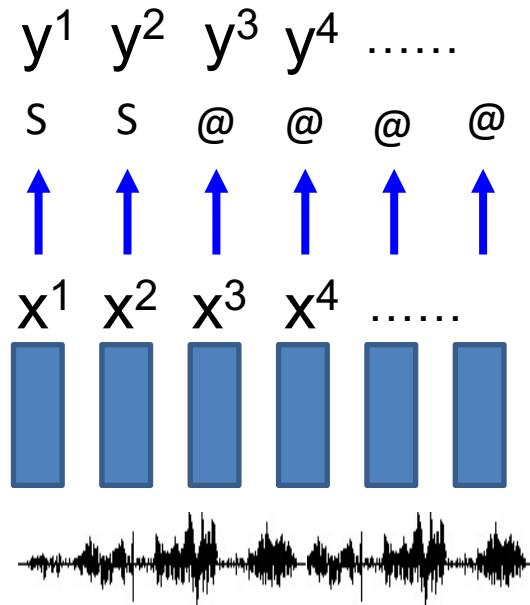
Application

- (Simplified) Speech Recognition



Utterance

1



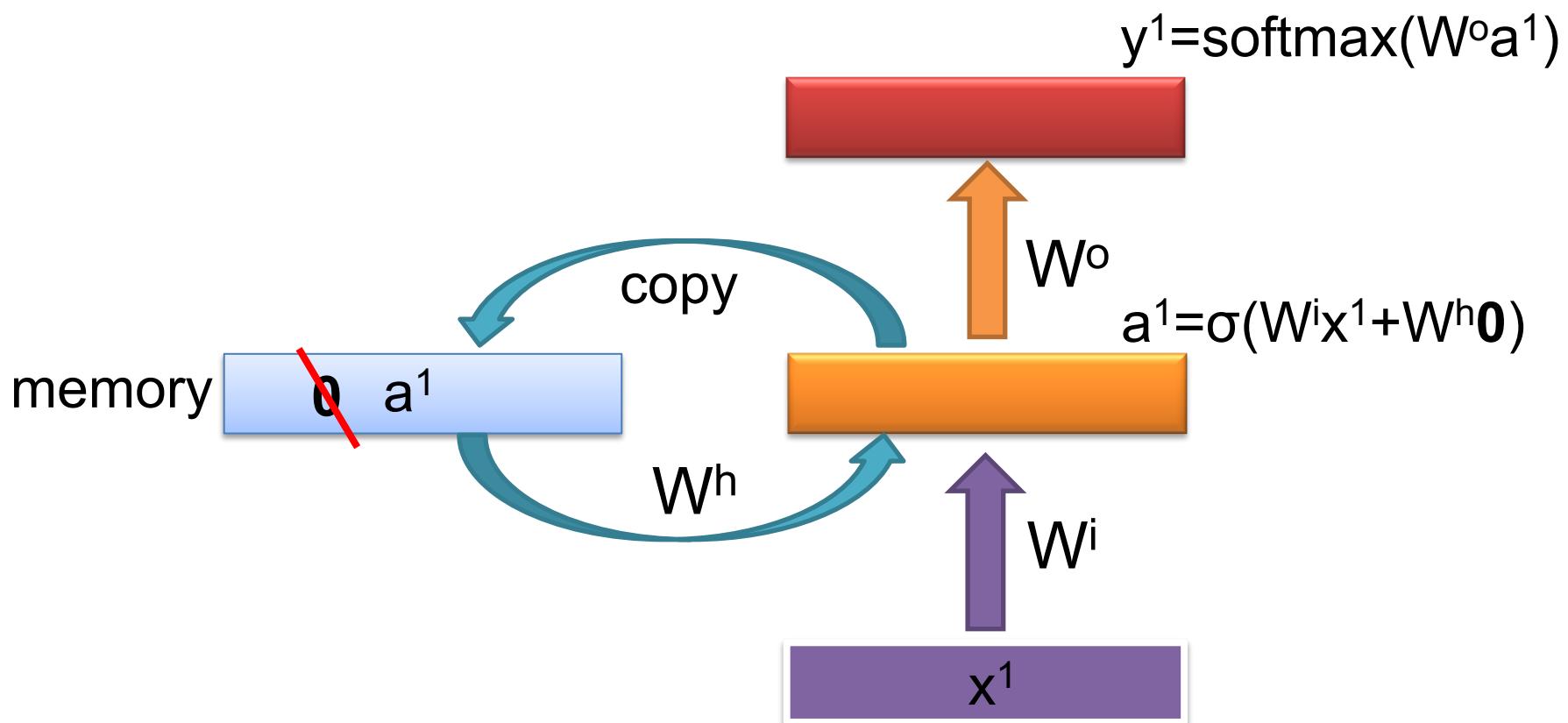
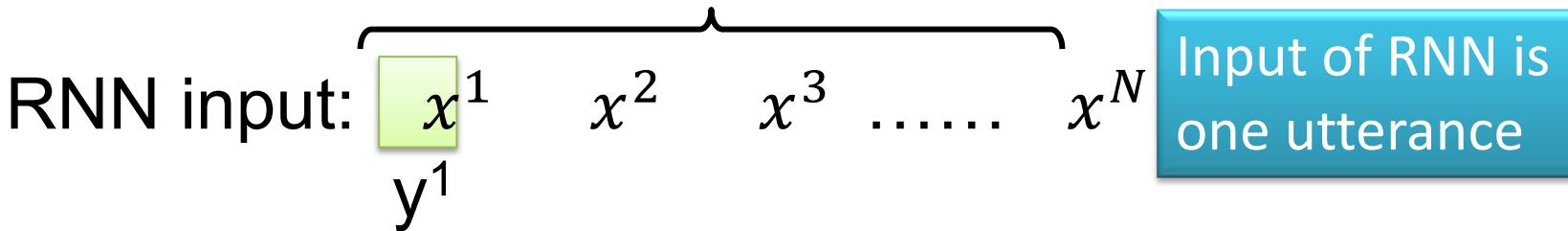
Utterance

2

We use DNN. All the frames are considered independently.

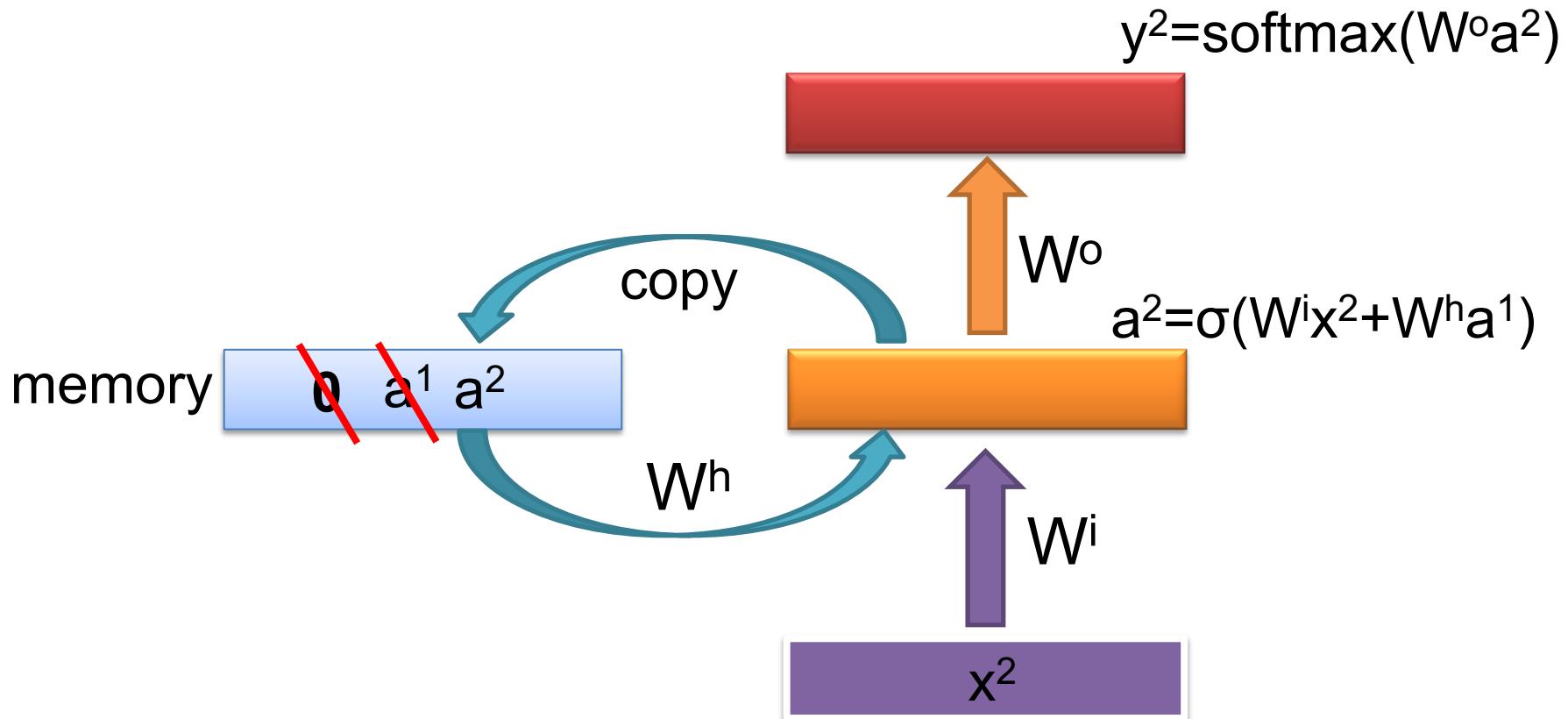
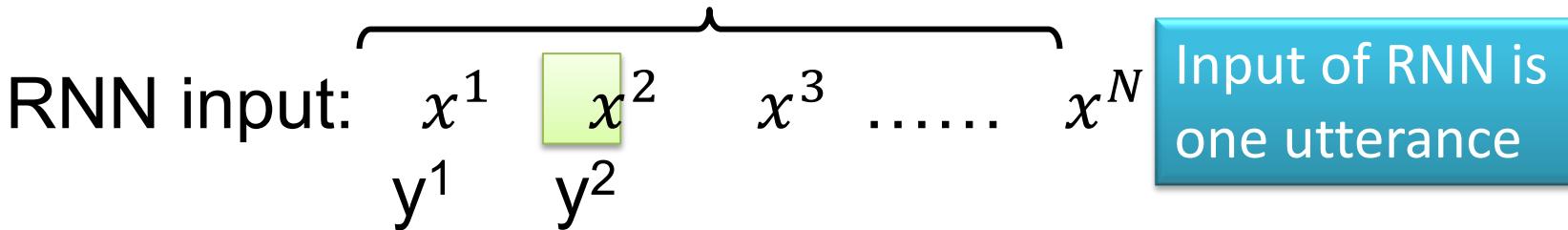
RNN

The order cannot change.



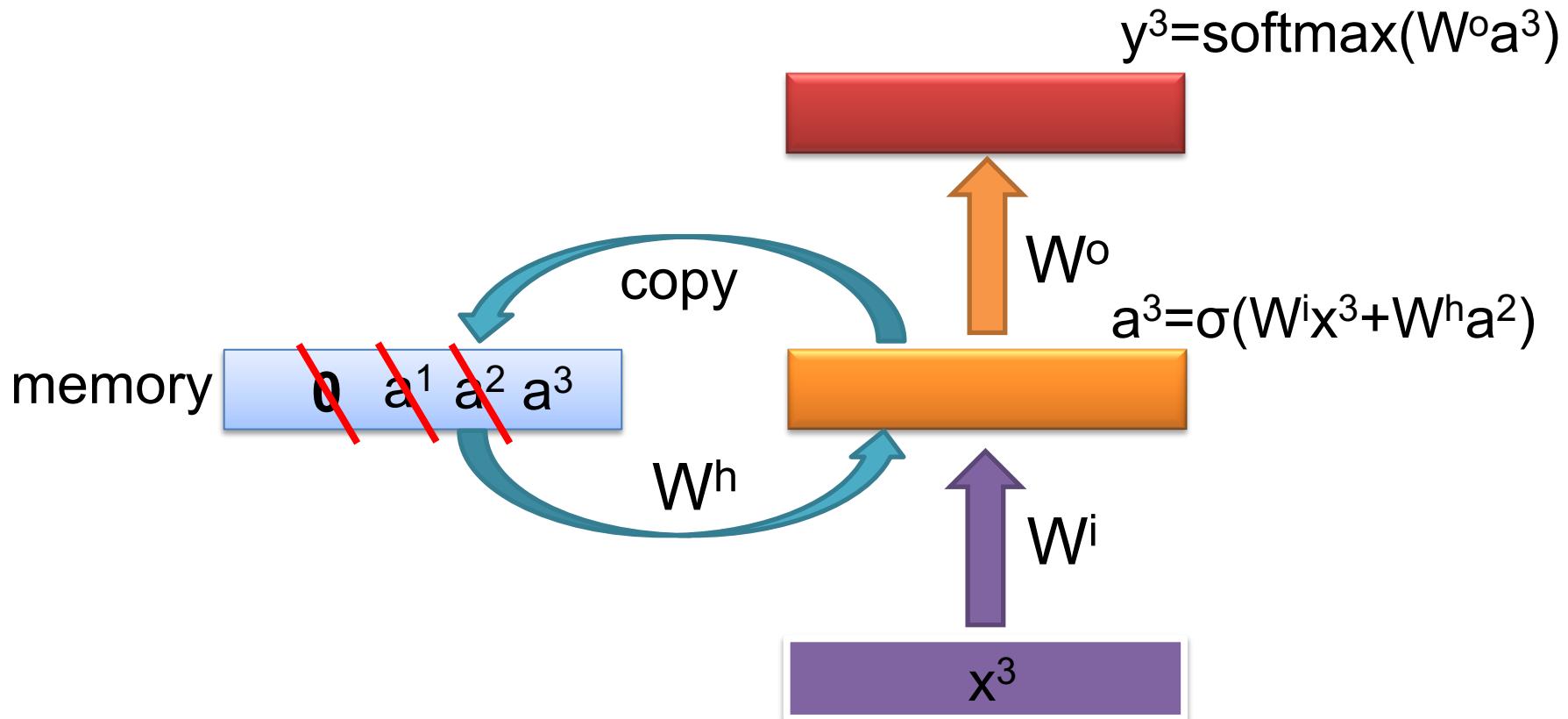
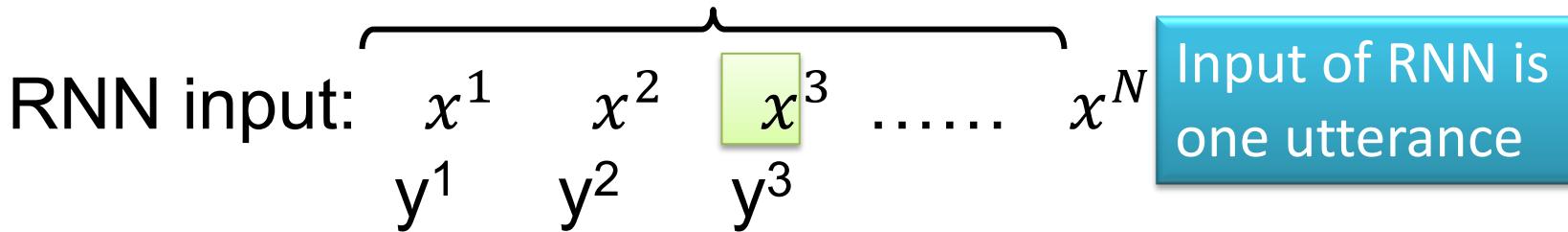
RNN

The order cannot change.



RNN

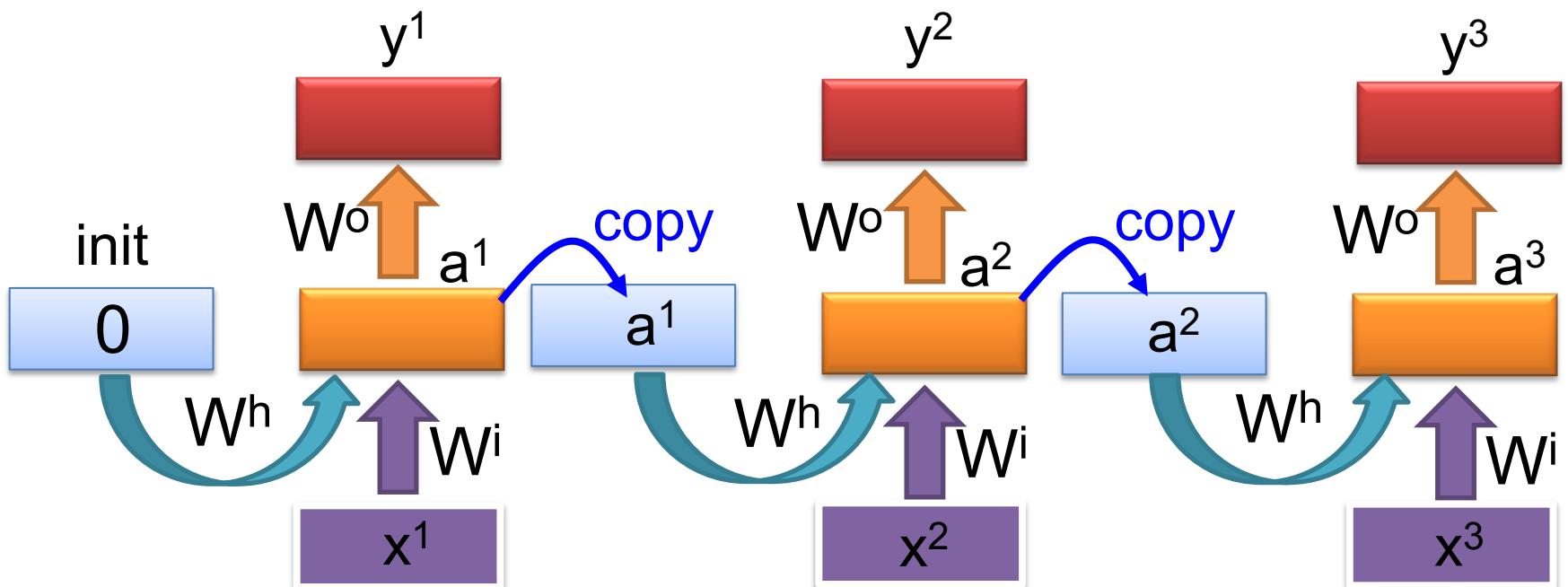
The order cannot change.



RNN

Input data: $x^1 \quad x^2 \quad x^3 \dots \quad x^N$

Input of RNN is one utterance



The same network is used again and again.

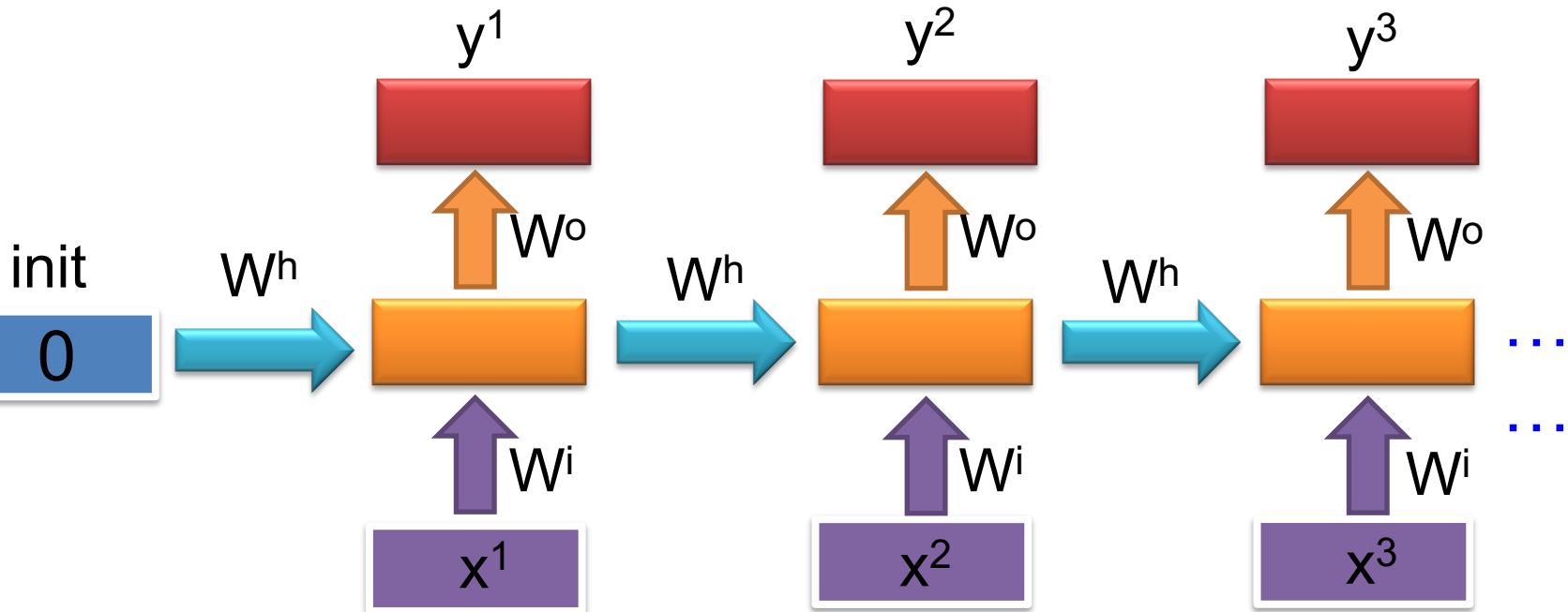
The values store in memory is different.

Output y^i depends on x^1, x^2, \dots, x^i

RNN

Input data: $x^1 \quad x^2 \quad x^3 \dots \dots \quad x^N$

Input of RNN is one utterance



The same network is used again and again.

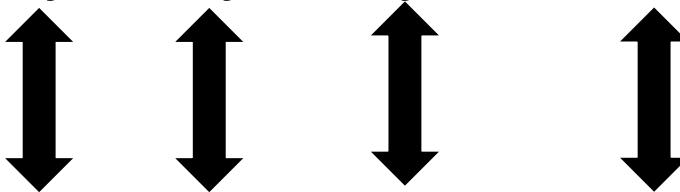
Output y^i depends on x^1, x^2, \dots, x^i



Cost

RNN input: $x^1 \quad x^2 \quad x^3 \dots \dots \quad x^N$

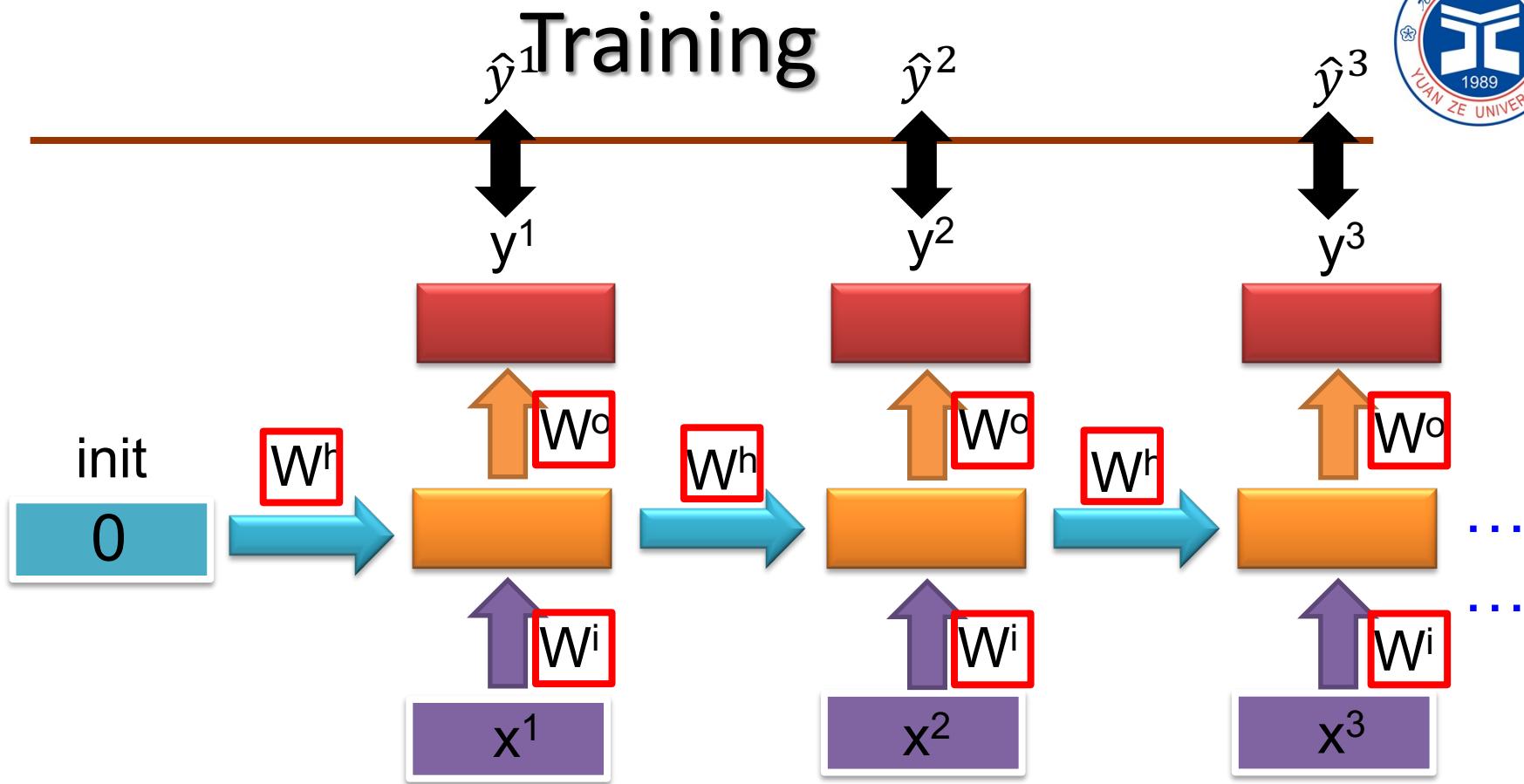
RNN output: $y^1 \quad y^2 \quad y^3 \dots \dots \quad y^N$



RNN output: $\hat{y}^1 \quad \hat{y}^2 \quad \hat{y}^3 \dots \dots \quad \hat{y}^N$

$$C = \frac{1}{2} \sum_{n=1}^N \|y^n - \hat{y}^n\|^2$$

$$C = \frac{1}{2} \sum_{n=1}^N -\log y_r^n$$



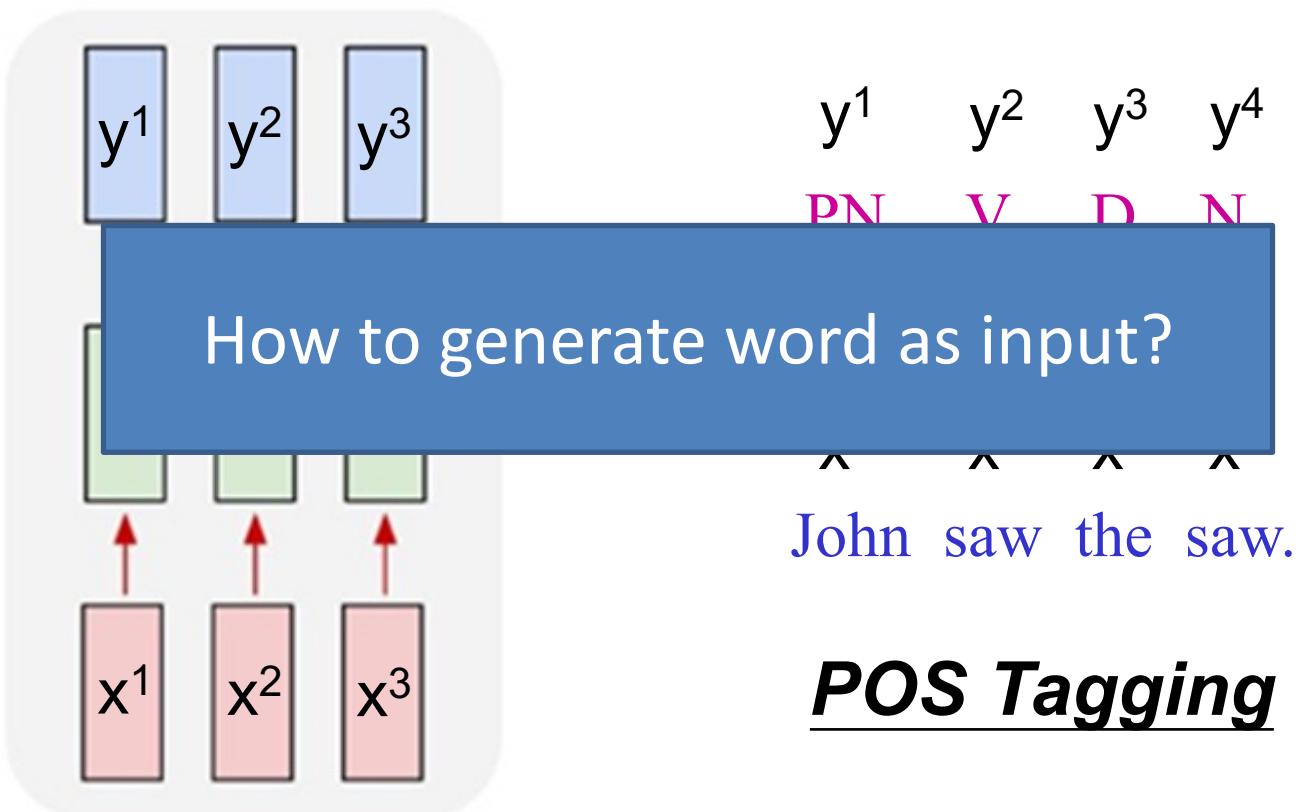
w is an element in W^h , W^i or W^o $\rightarrow w \leftarrow w - \eta \partial C / \partial w$

→ Backpropagation through time (BPTT)

RNN Training is very difficult in practice.

Applications Suitable for RNN

- Input and output are vector sequences with *the same length*





More Applications

- Name entity recognition
 - Identifying names of people, places, organizations, etc. from a sentence
 - Harry Potter is a student of Hogwarts and lived on Privet Drive.
 - people, organizations, places, not a name entity
- Information extraction
 - Extract pieces of information relevant to a specific application, e.g. flight booking
 - I would like to leave Boston on November 2nd and arrive in Taipei before 2 p.m.
 - place of departure, destination, time of departure, time of arrival, other



Word Representation:

1-of-N Encoding

- We let machine read a lot of articles and use one hot encoding on each word.

"a"	"abbreviations"	"zoology"	"zoom"
1	0	0	0
0	1	0	1
0	0	0	0
.	.	.	.
.	.	.	.
.	.	.	.
0	0	0	0
0	0	1	0
0	0	0	0
0	0	0	1



Drawback of 1-of-N Encoding

- Dog is not close to cat
 - Machine can not understand the meaning of word.
- Waste a lot of entity
 - Most of entity are zero
 - Dimension of vector = vocabulary size (very large)

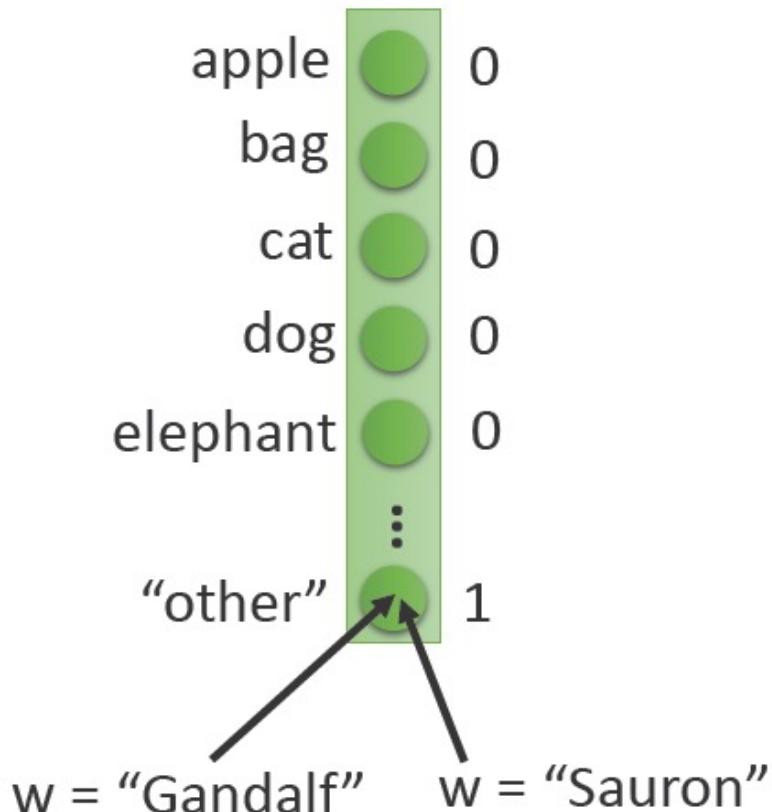


Drawback of 1-of-N Encoding

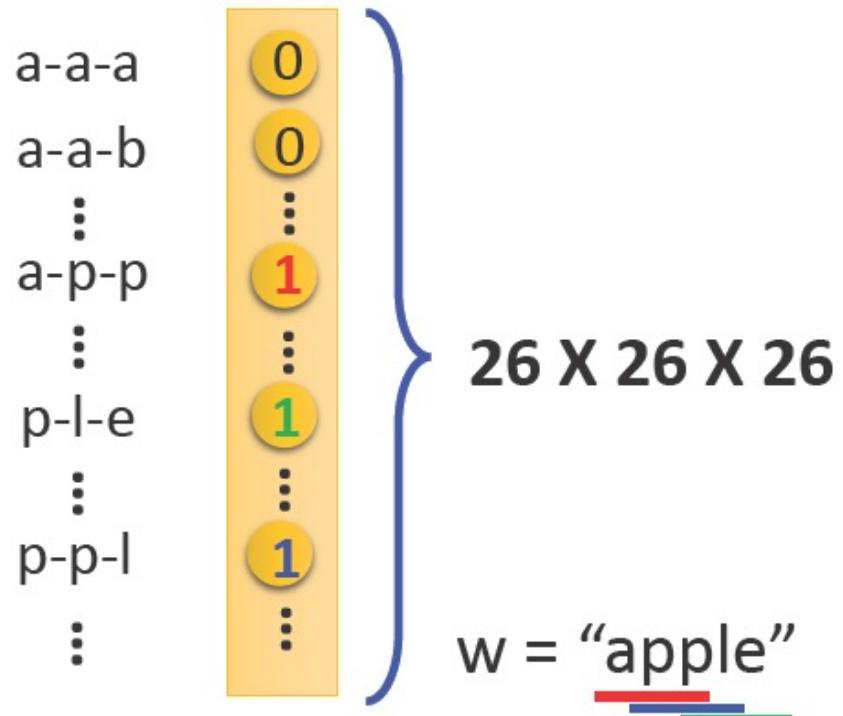
- Dog is not close to cat
 - Machine can not understand the meaning of word.
- Waste a lot of entity
 - Most of entity are zero
 - Dimension of vector = vocabulary size (very large)

Beyond 1-of-N Encoding

Dimension for “Other”



Word hashing





Outline

Recurrent Neural Network (RNN)



Variants of RNN

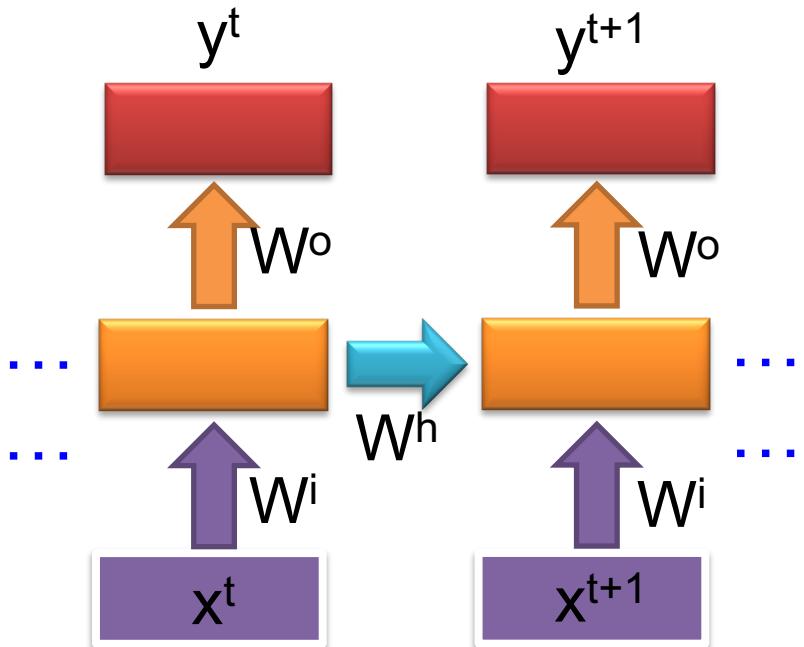


Long Short-term Memory (LSTM)

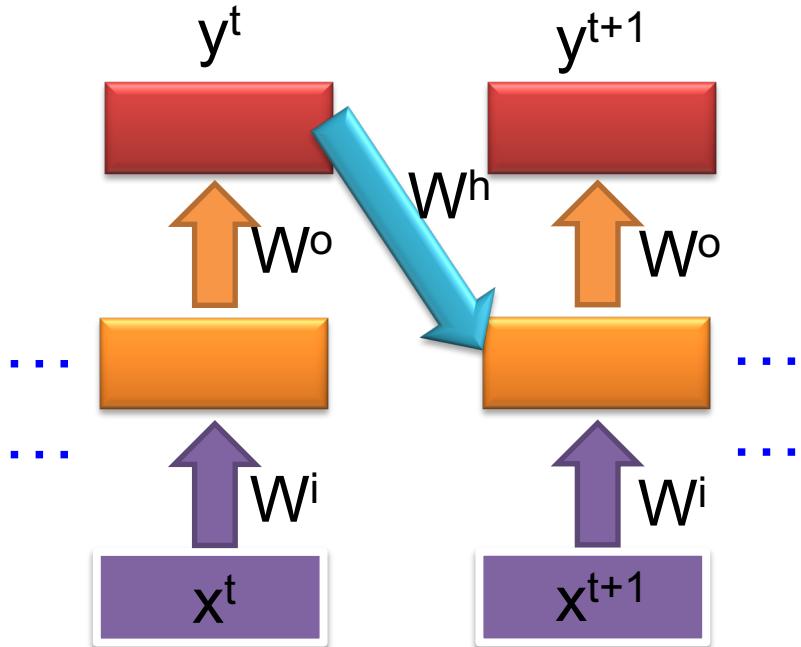
Elman Network & Jordan Network



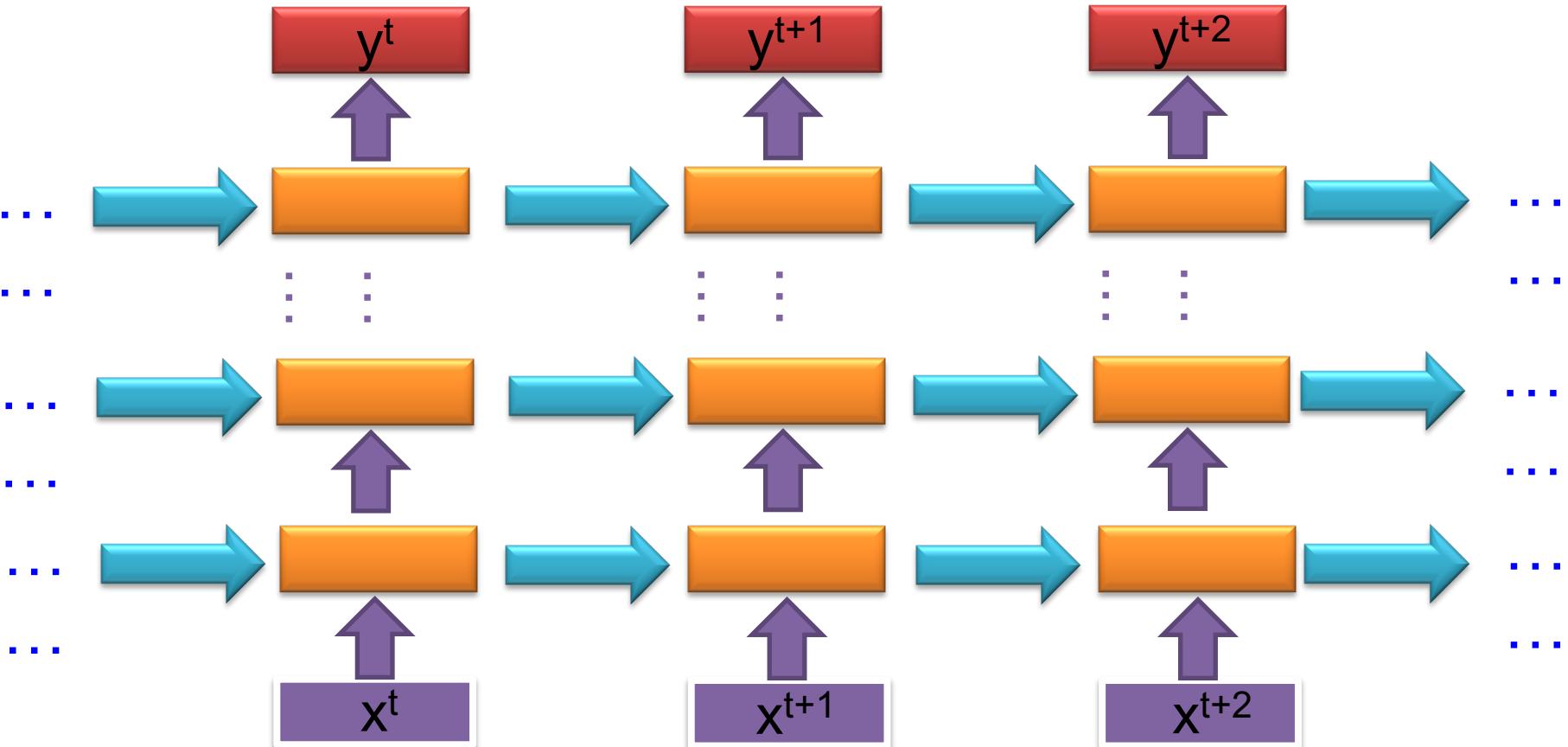
Elman Network



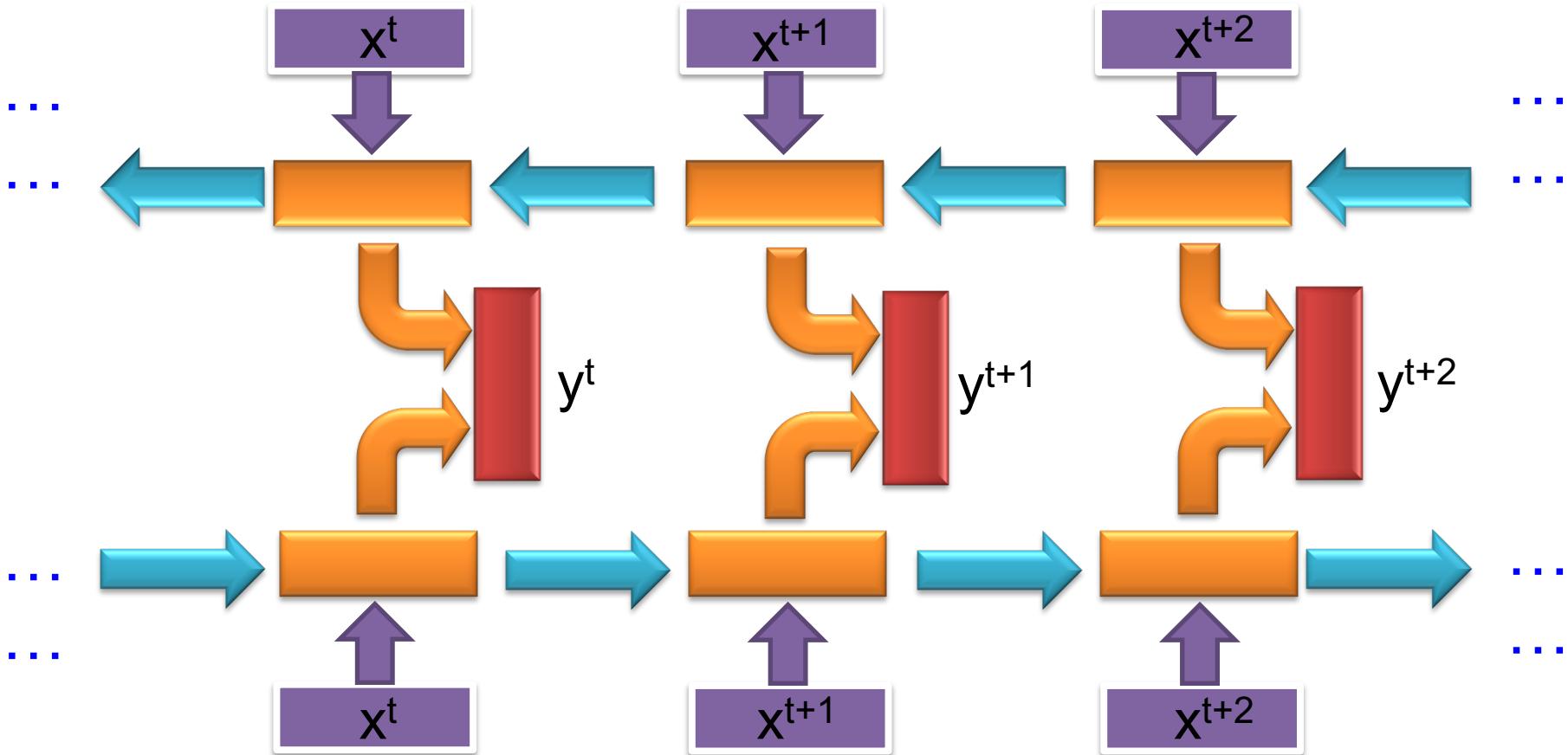
Jordan Network



Deep RNN



Bidirectional RNN





Outline

Recurrent Neural Network (RNN)

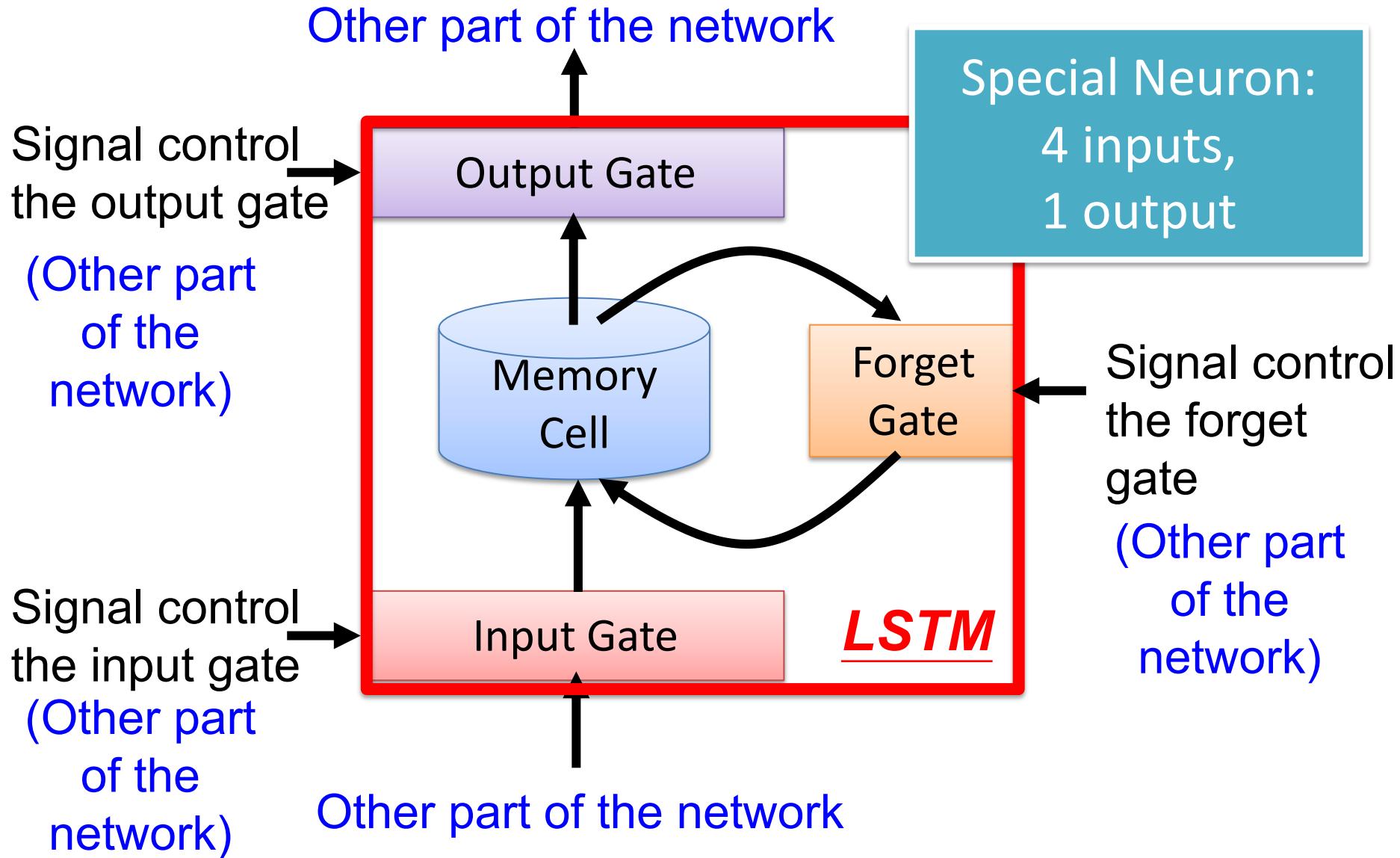


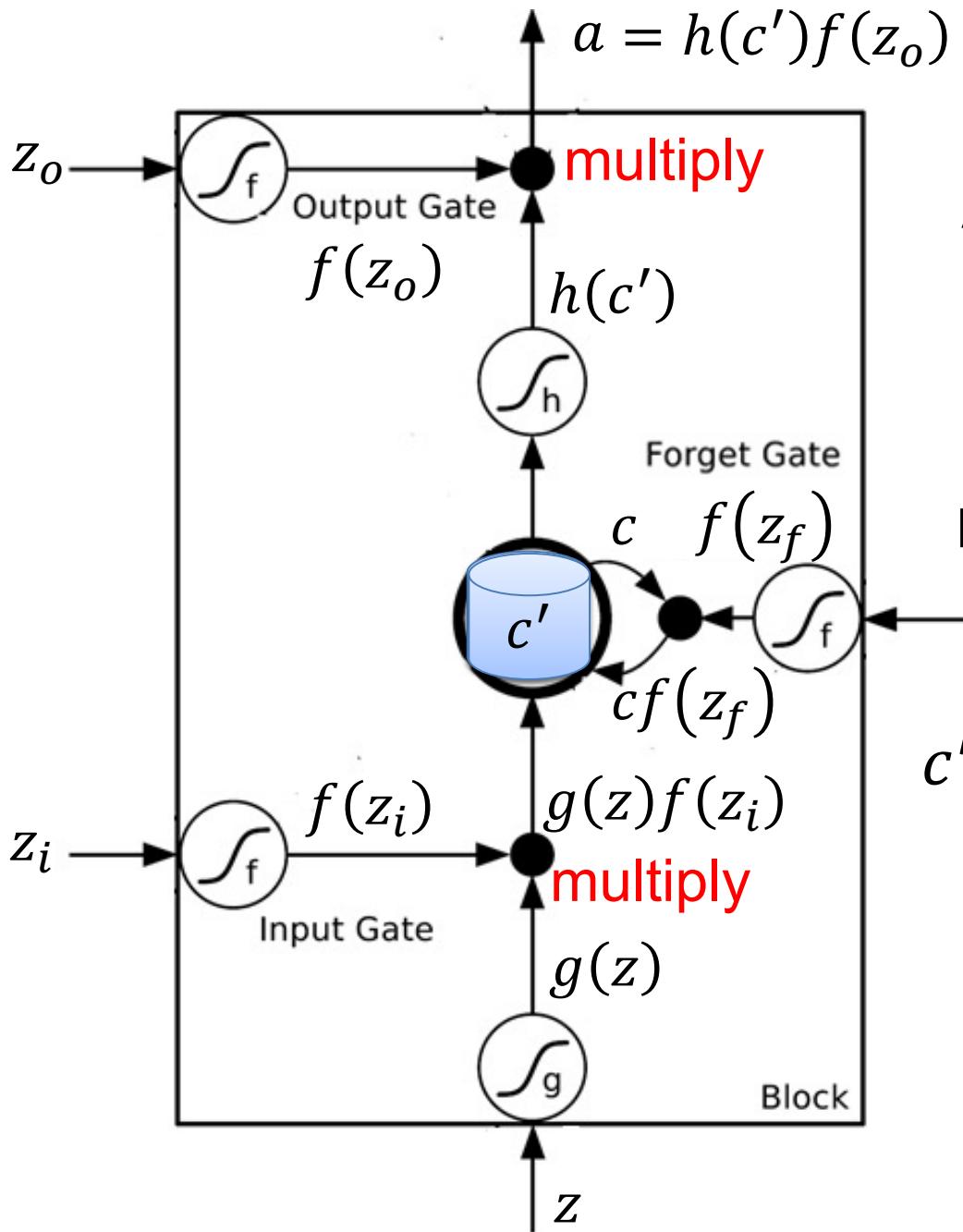
Variants of RNN



Long Short-term Memory (LSTM)

Long Short-term Memory (LSTM)





Activation function f is usually a sigmoid function

Sigmoid function between 0 and 1

Mimic open and close gate

$$z_f$$

$$c' = g(z)f(z_i) + cf(z_f)$$

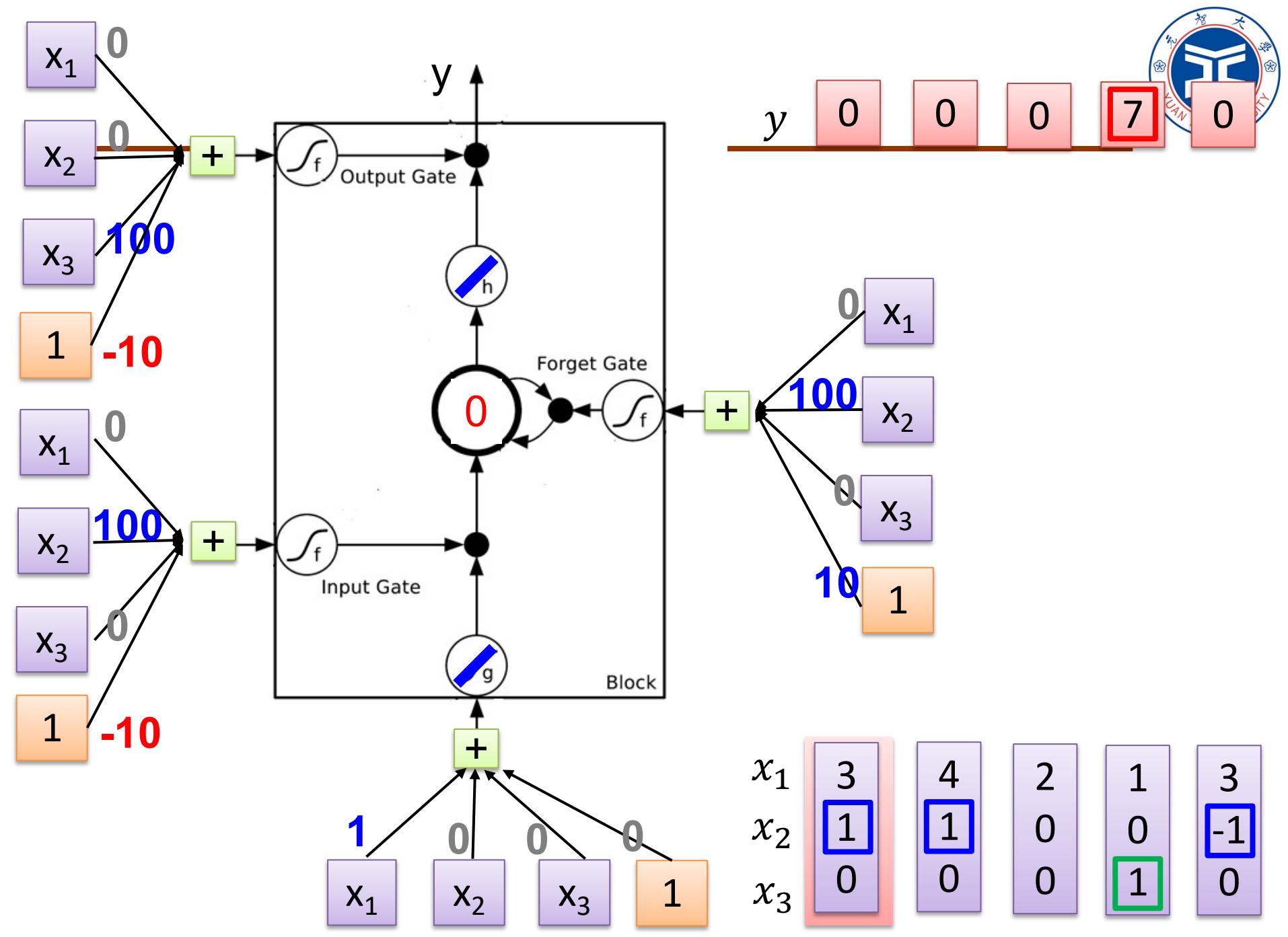
LSTM - Example

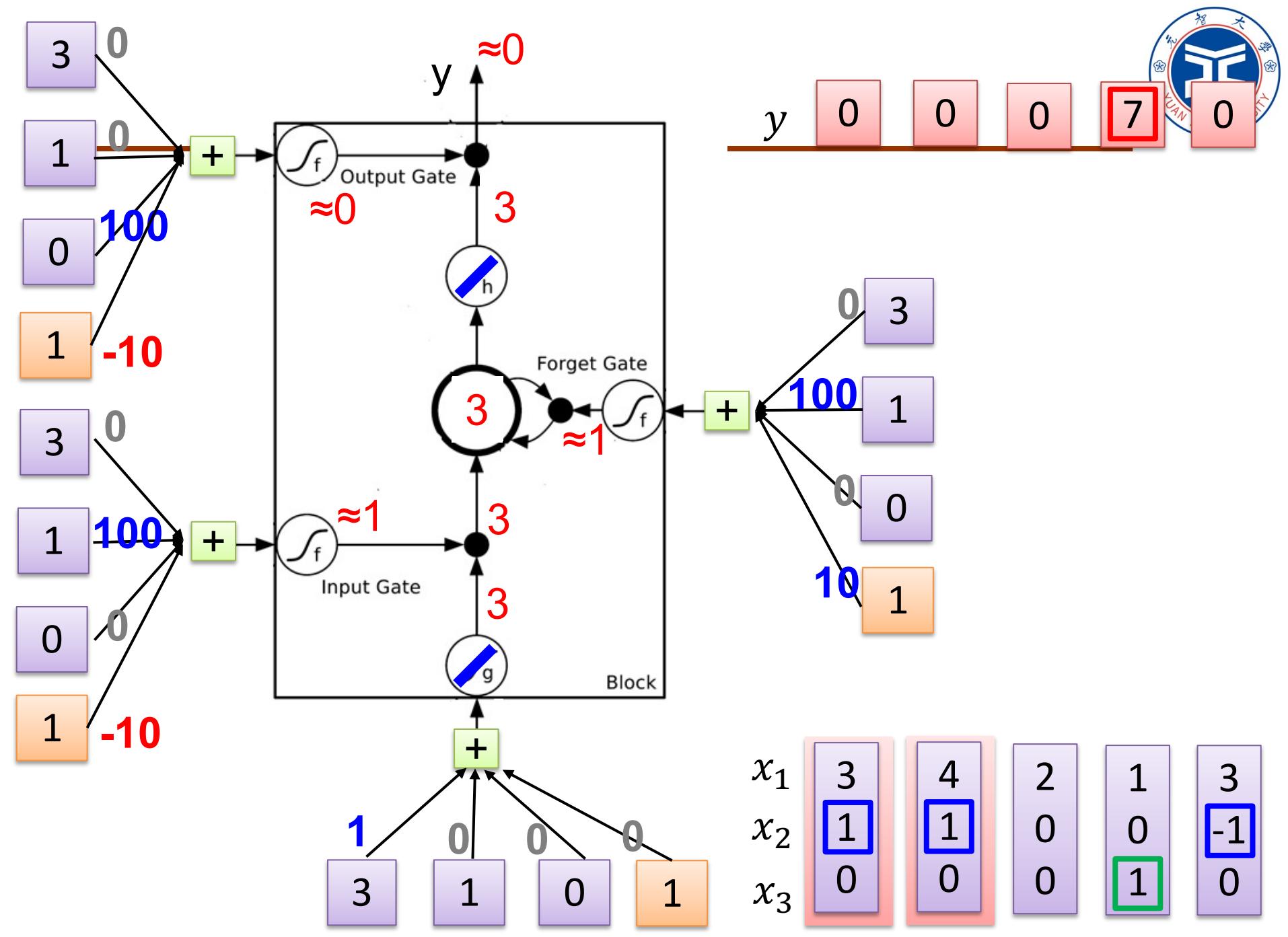
	0	0	3	3	7	7	7	0	6
x_1	1	3	2	4	2	1	3	6	1
x_2	0	1	0	1	0	0	-1	1	0
x_3	0	0	0	0	0	1	0	0	1
y	0	0	0	0	0	7	0	0	6

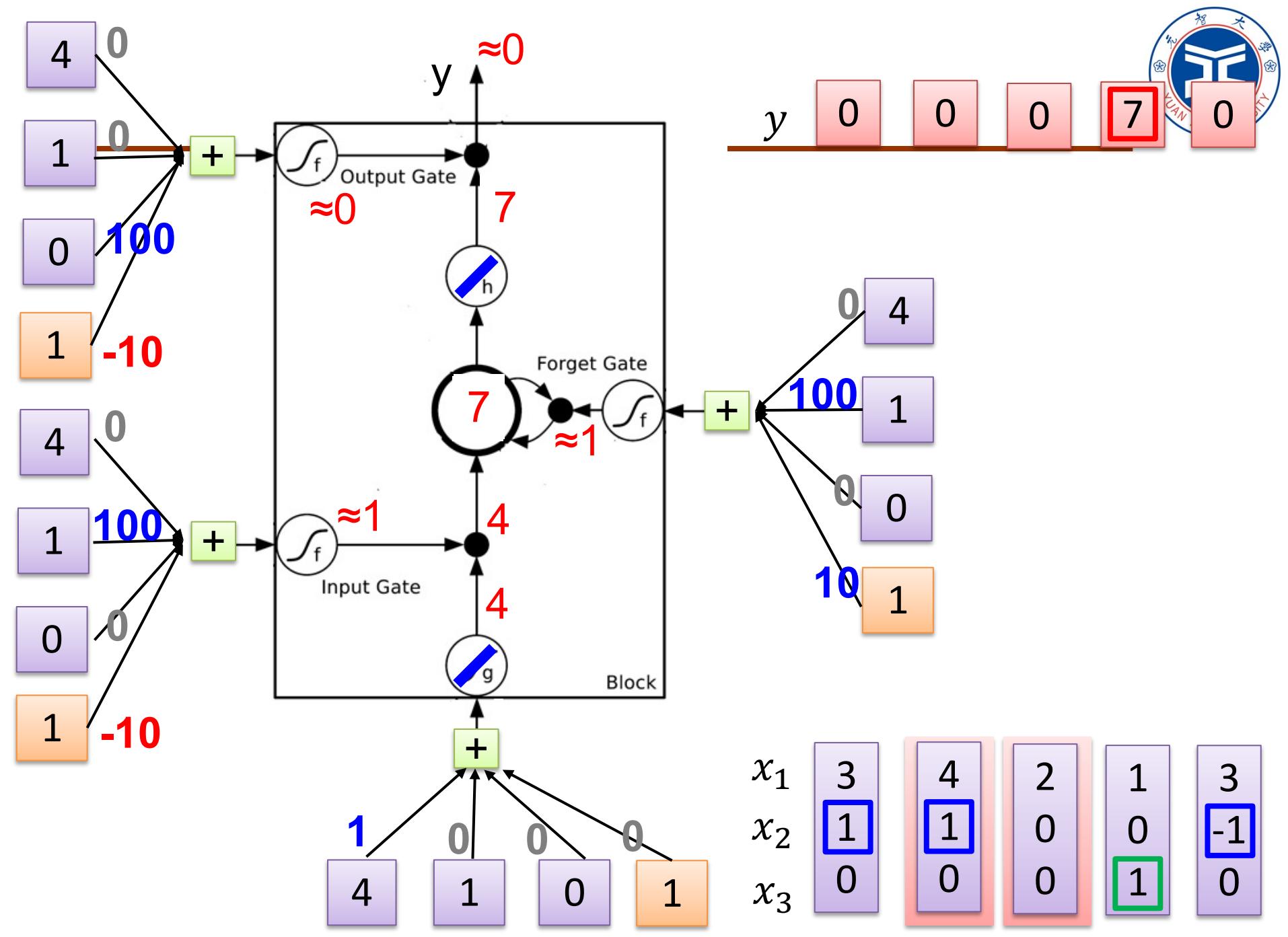
When $x_2 = 1$, add the numbers of x_1 into the memory

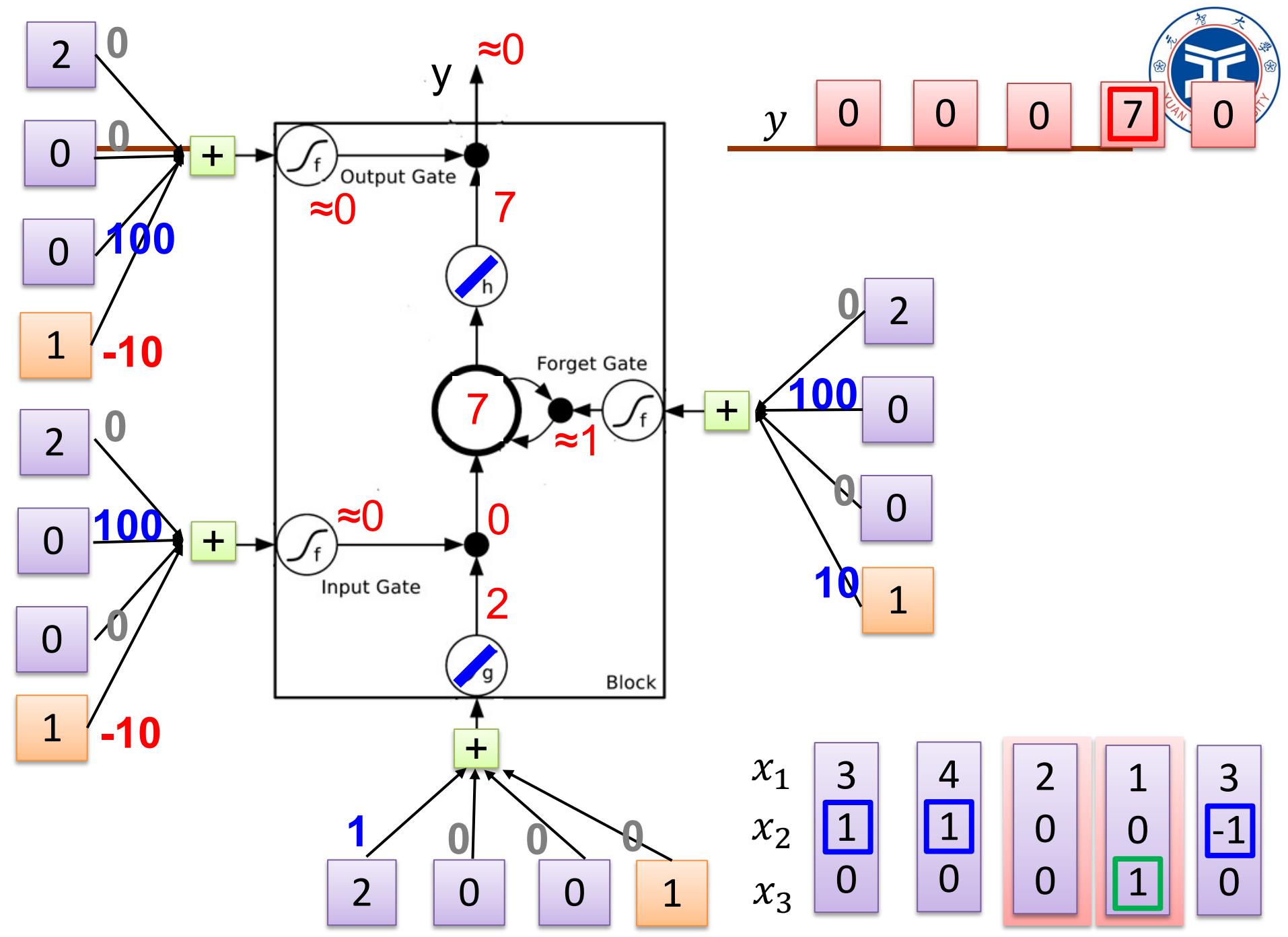
When $x_2 = -1$, reset the memory

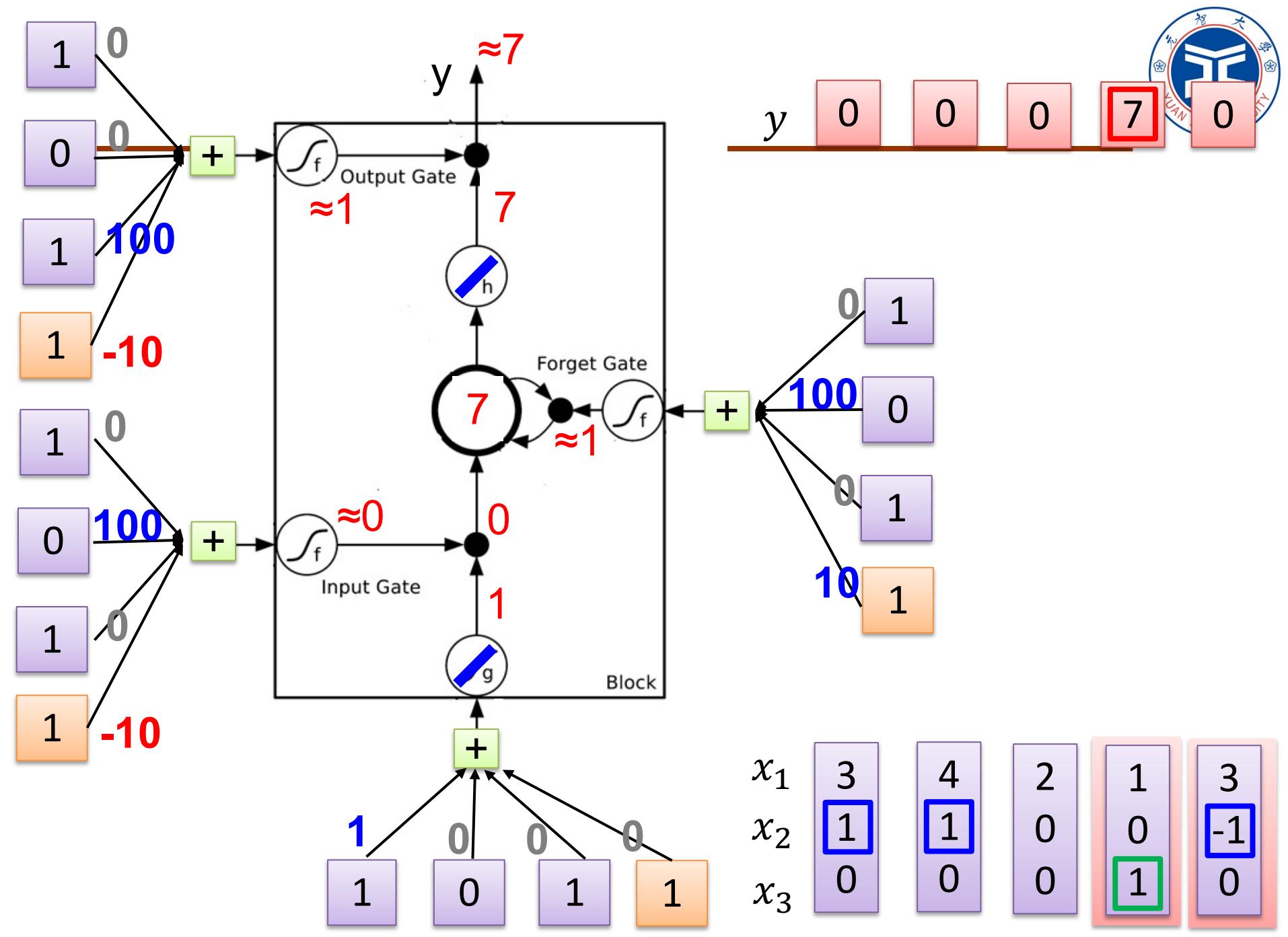
When $x_3 = 1$, output the number in the memory.

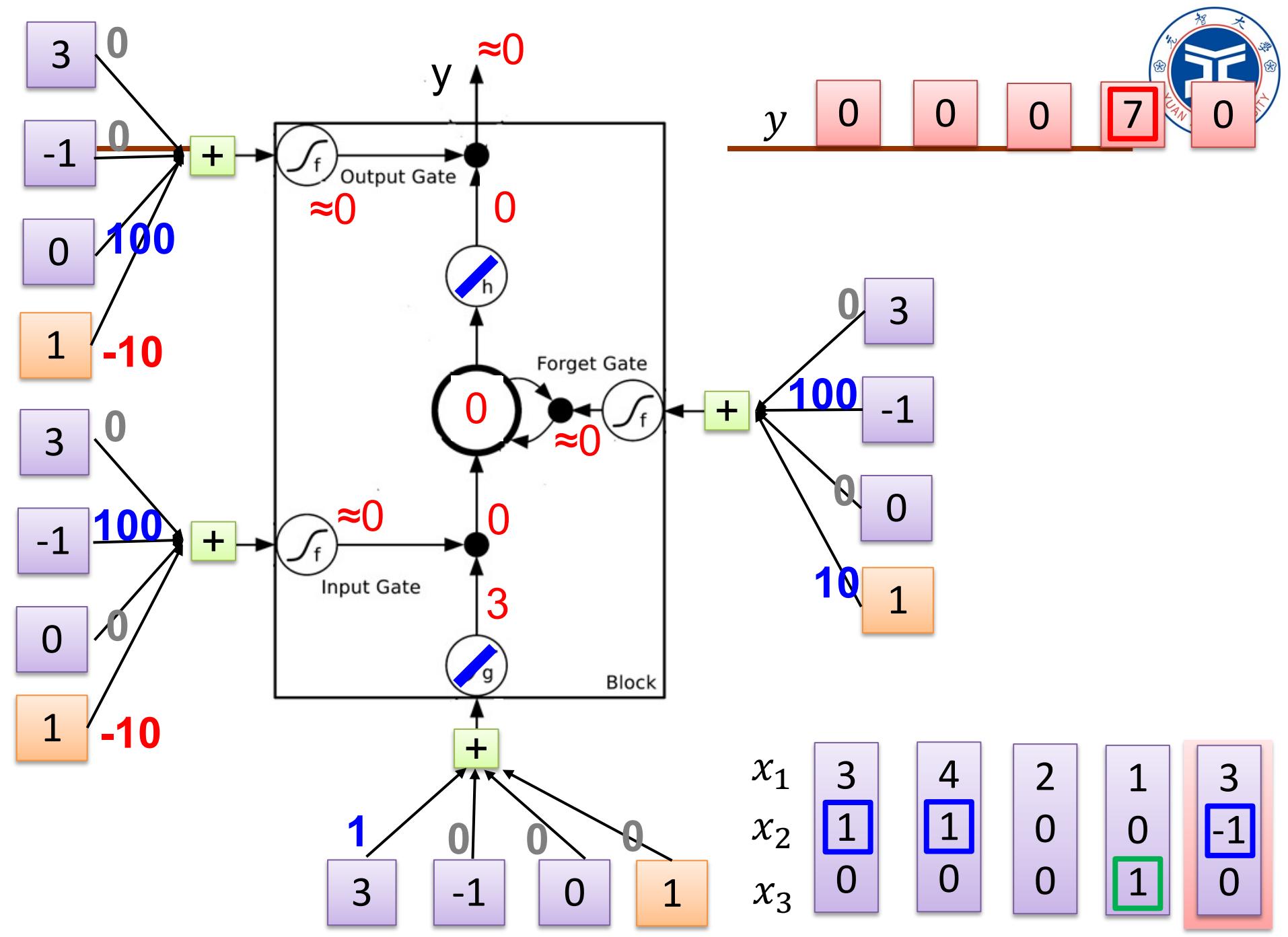








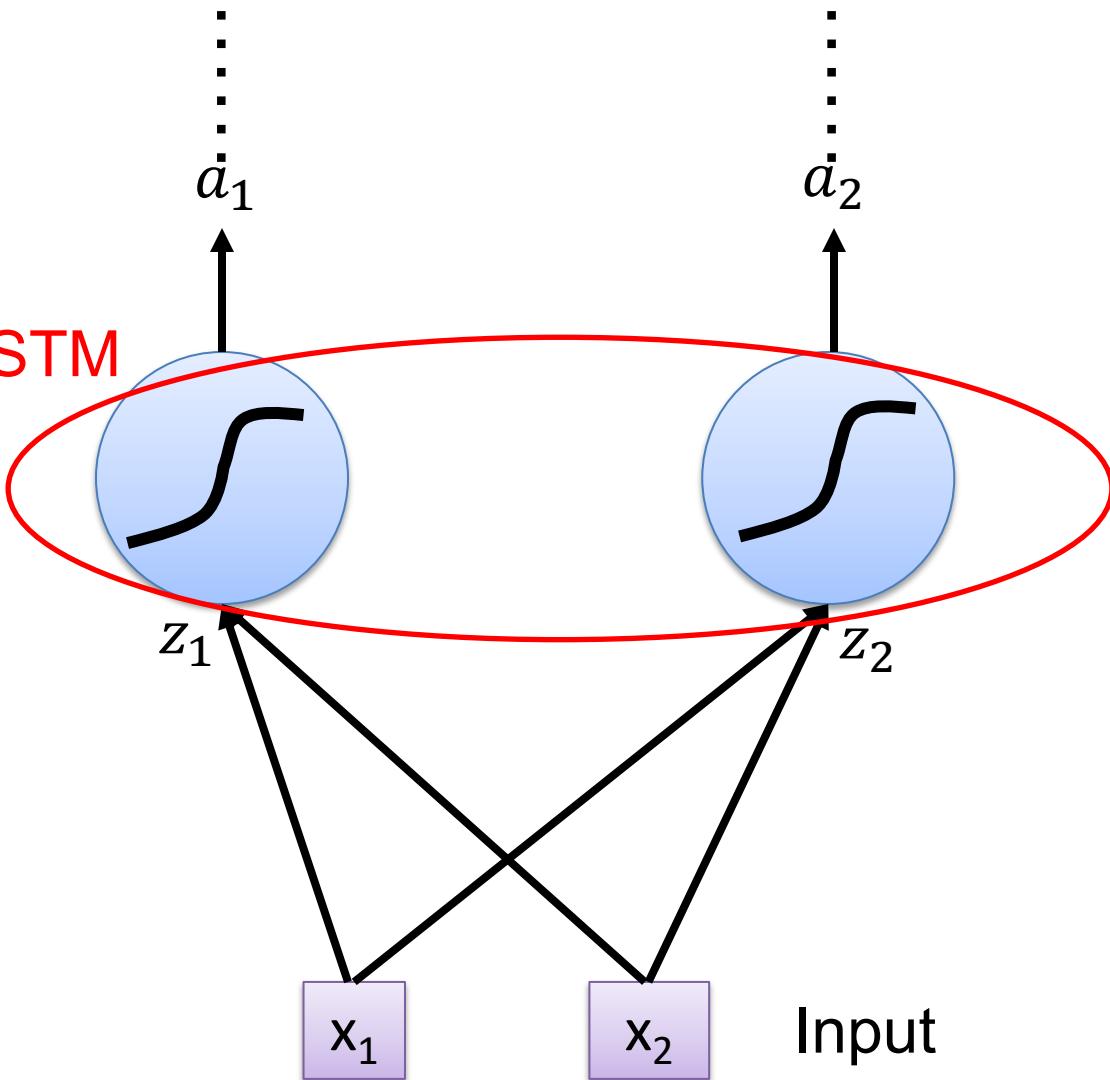


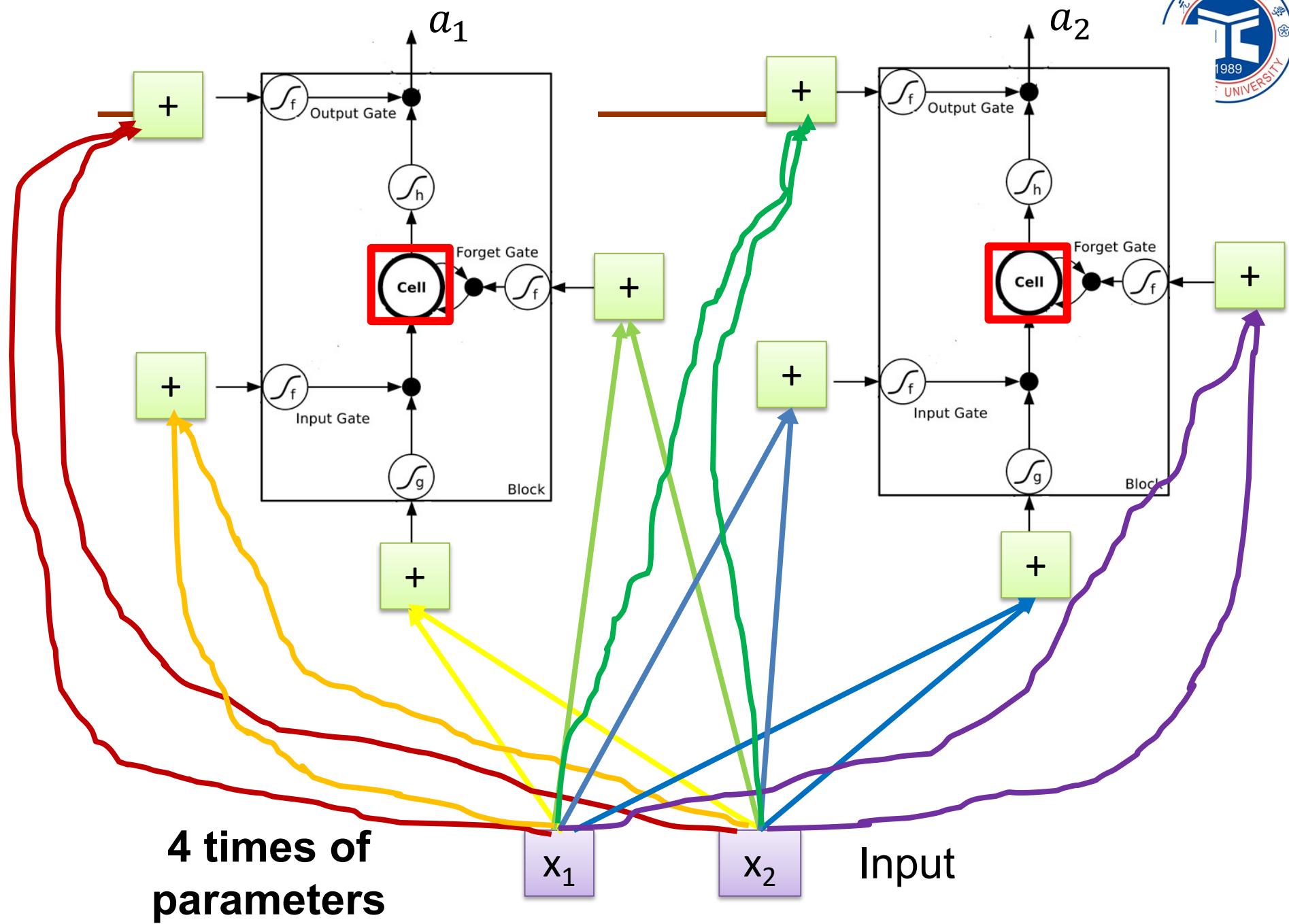


Original Network

- Simply replace the neurons with LSTM

Replace with LSTM

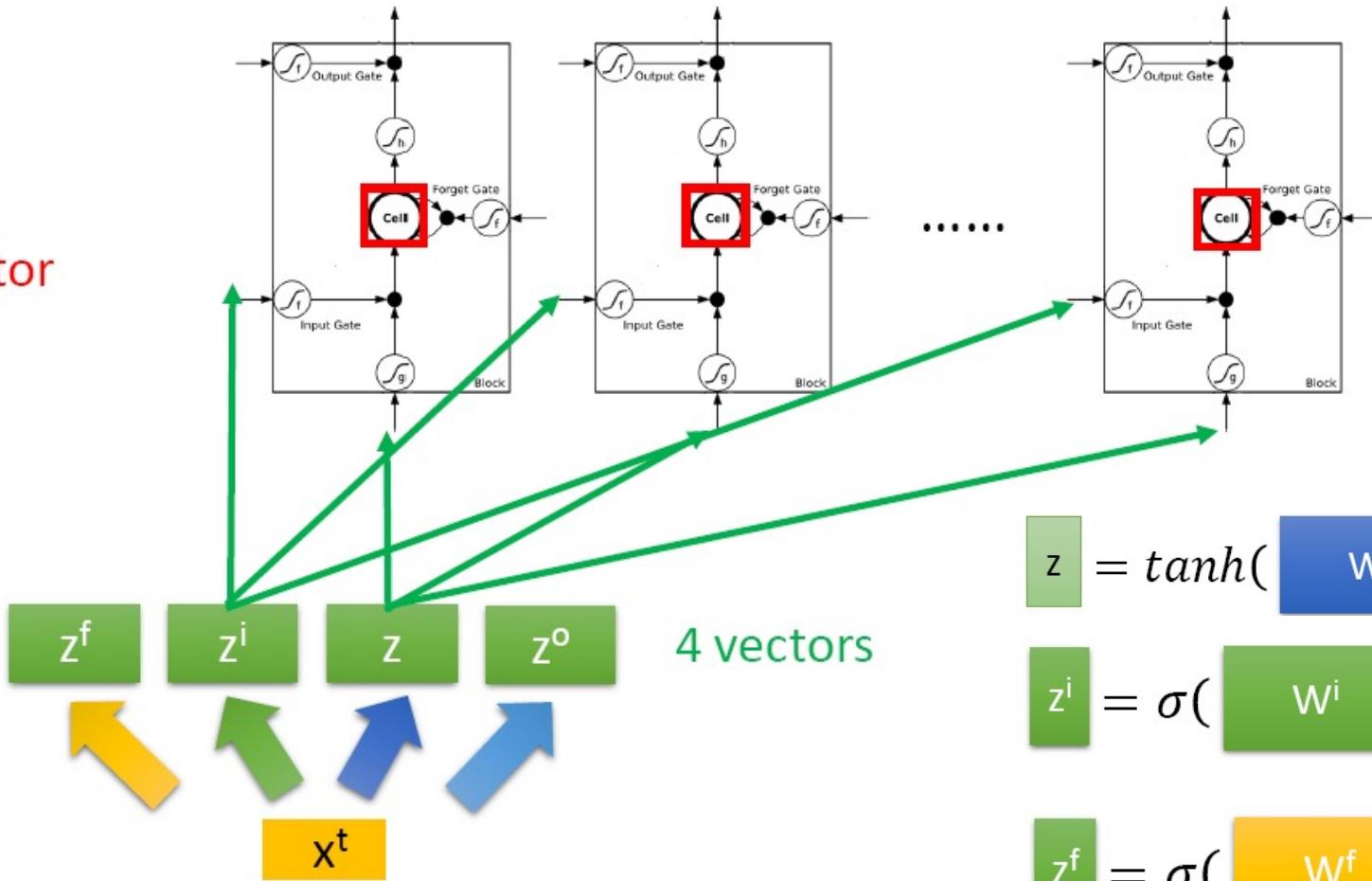




LSTM

c^{t-1}

vector



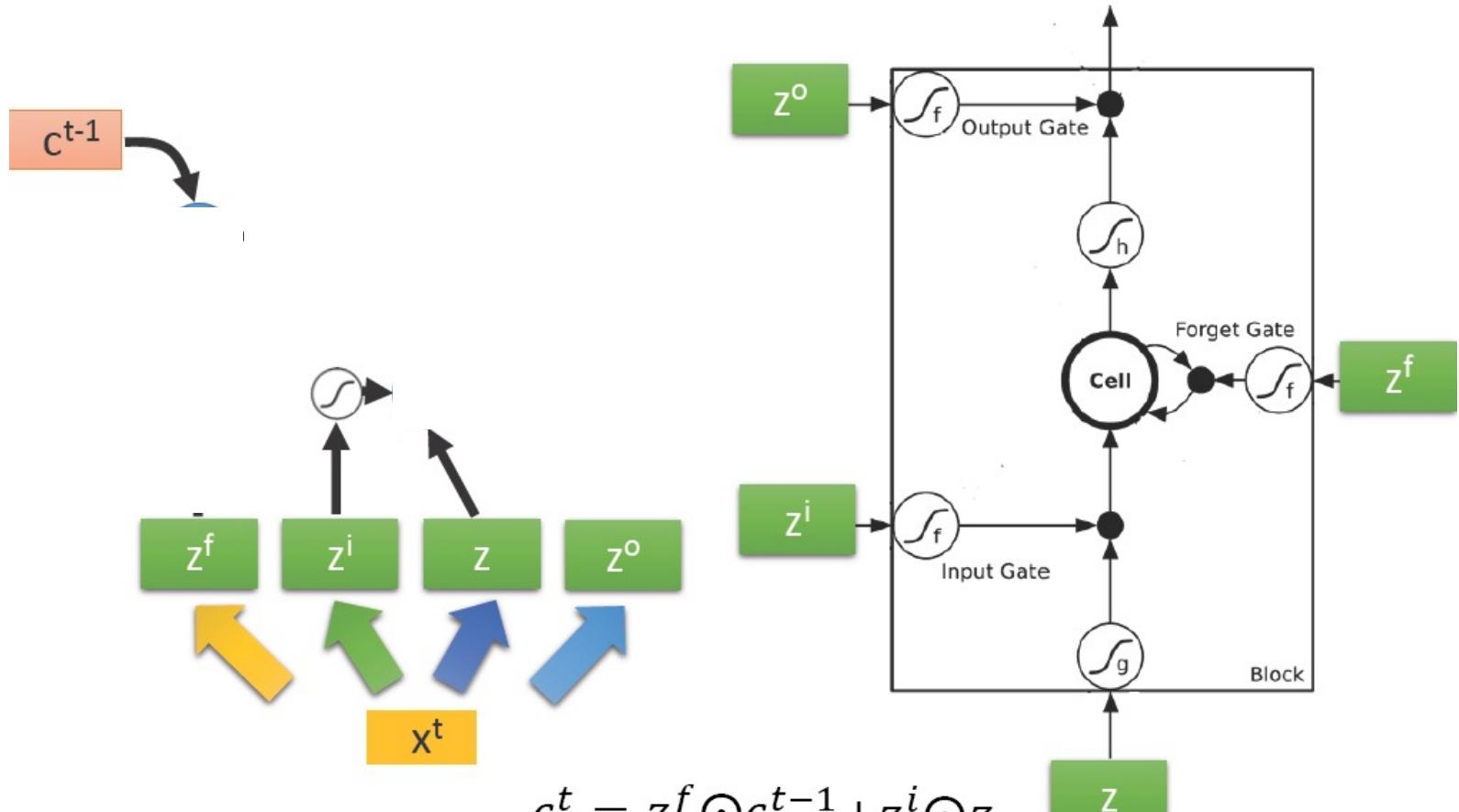
$$z = \tanh(\begin{matrix} W \\ x^t \end{matrix})$$

$$z^i = \sigma(\begin{matrix} W^i \\ x^t \end{matrix})$$

$$z^f = \sigma(\begin{matrix} W^f \\ x^t \end{matrix})$$

$$z^o = \sigma(\begin{matrix} W^o \\ x^t \end{matrix})$$

LSTM



$$c^t = z^f \odot c^{t-1} + z^i \odot z$$

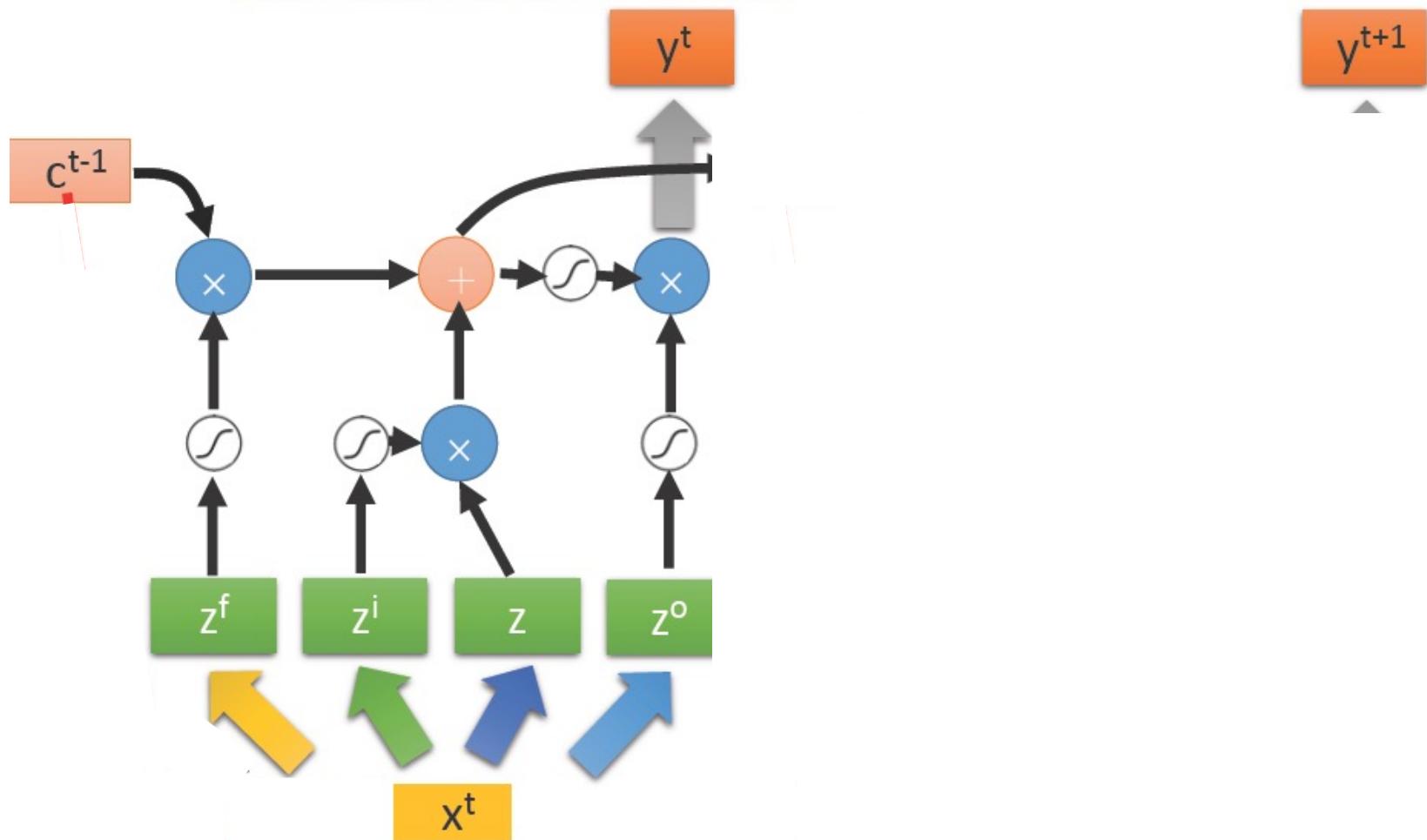
$$h^t = z^o \odot \tanh(c^t)$$

$$y^t = \sigma(W' h^t)$$

LSTM

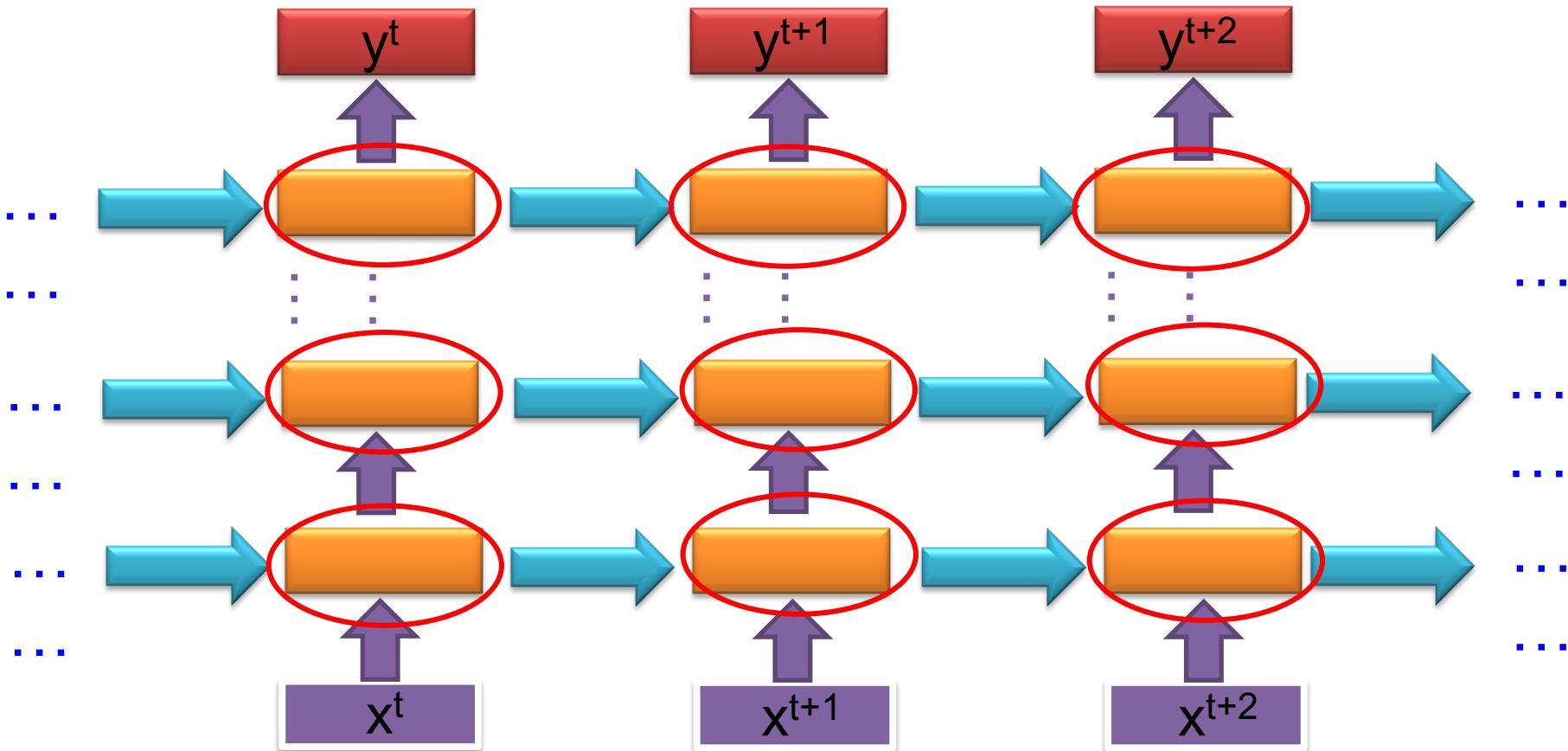
Complex LSTM add last output as next input!

Peephole: last memory as next input!!

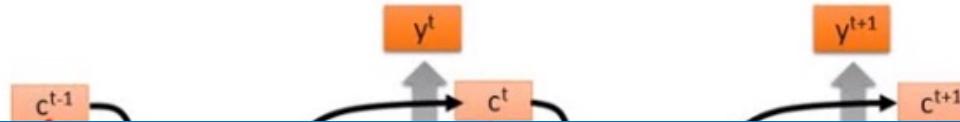


Deep RNN

- Replace neuron as LSTM



Multi-layer LSTM

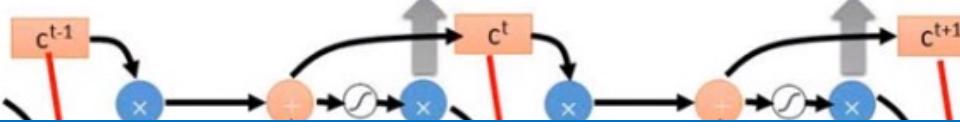


Too Complicated!!

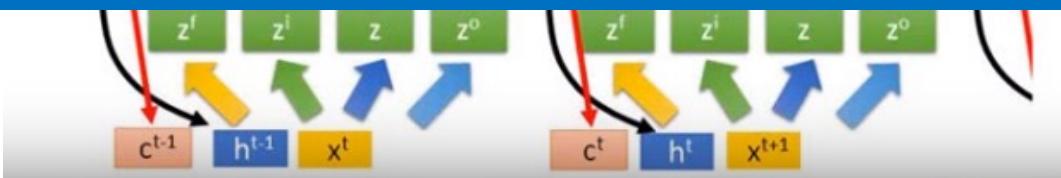
LSTM becomes a standard library.
LSTM almost can represent RNN.



Find “LSTM” in Keras.
You can import LSTM immediately.



Keras supports
“LSTM”, “GRU”, “SimpleRNN” layers

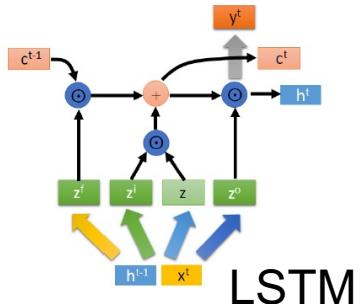


Keras is an [open-source neural-network](#) library written in [Python](#).

It is capable of running on top of [TensorFlow](#), [Microsoft Cognitive Toolkit](#), [Theano](#), or [PlaidML](#)

Gate Recurrent Unit (GRU)

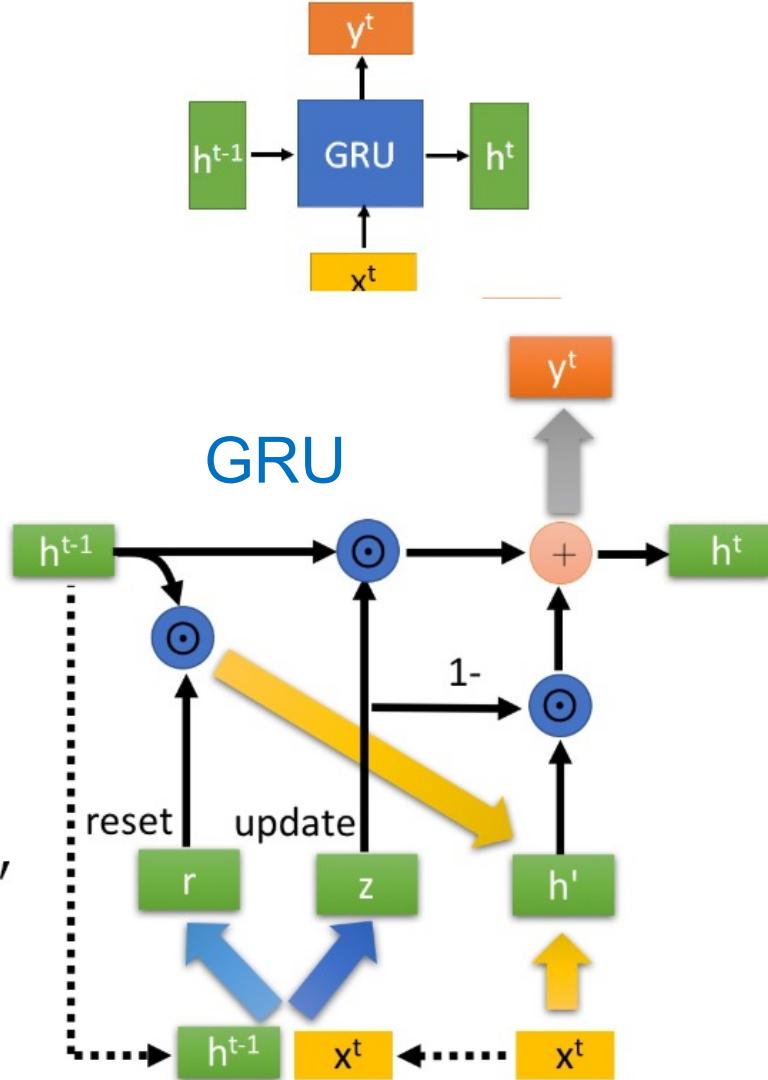
- Simplify LSTM
- Only two gate (reset, update gate)
 - parameters are fewer than LSTM
 - Prevent overfitting
- Input gate and forget gate are linked.
 - When there is new input, forget the older value.



$$h^t = z \odot h^{t-1} + (1 - z) \odot h'$$

Similar

$$c^t = z^f \odot c^{t-1} + z^i \odot z$$

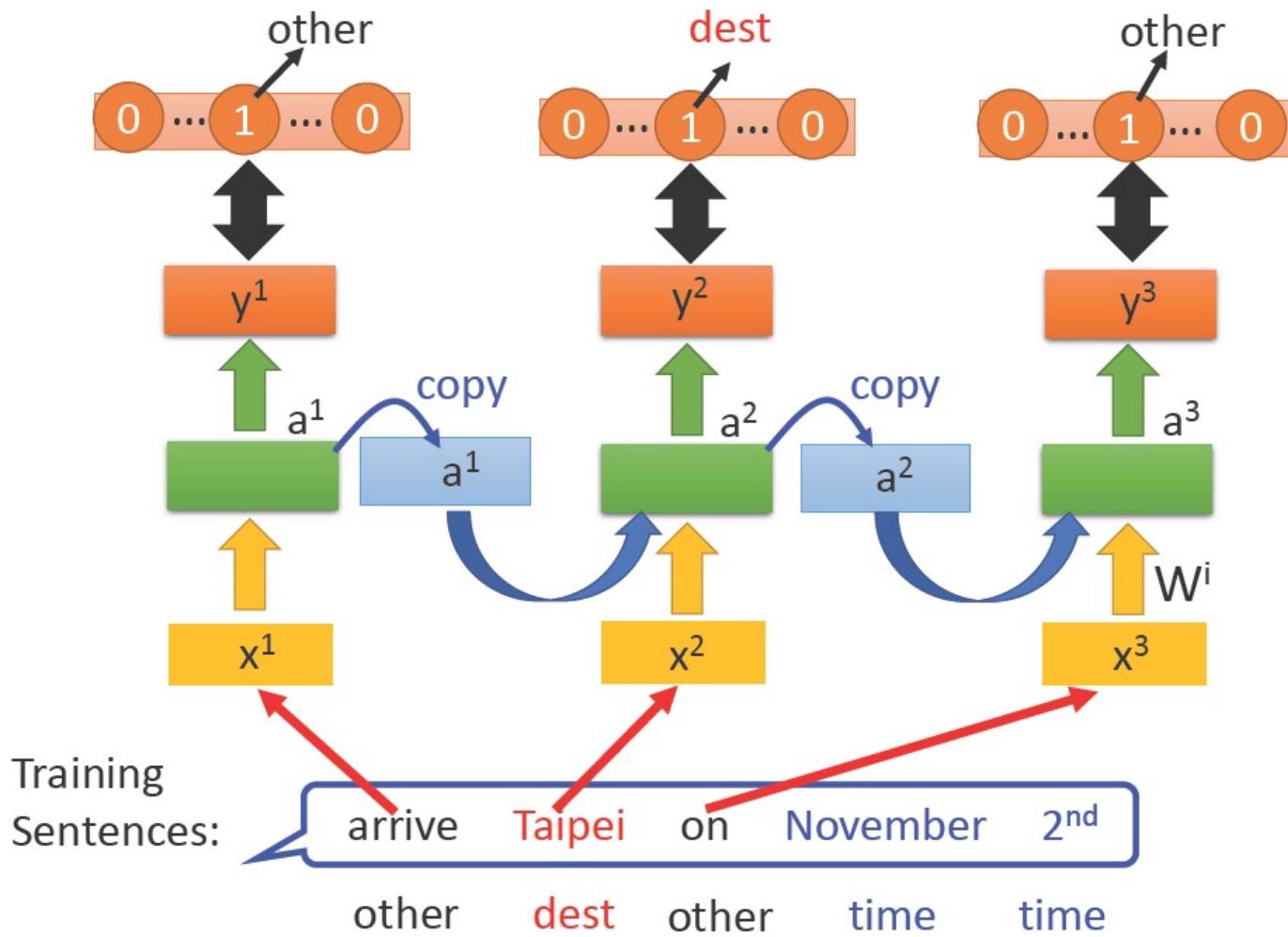




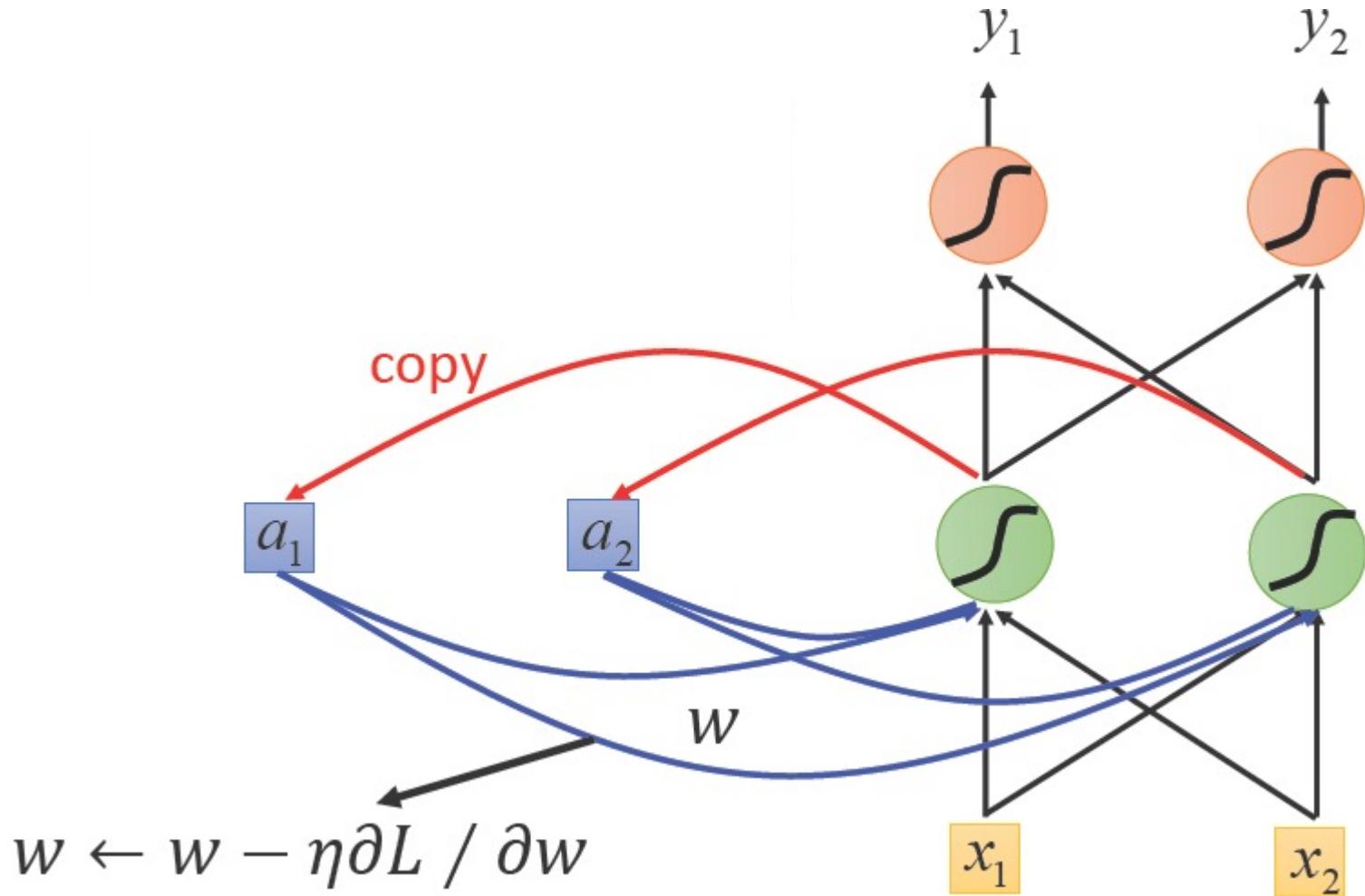
Learning in RNN

Learning Target

- Slot filling

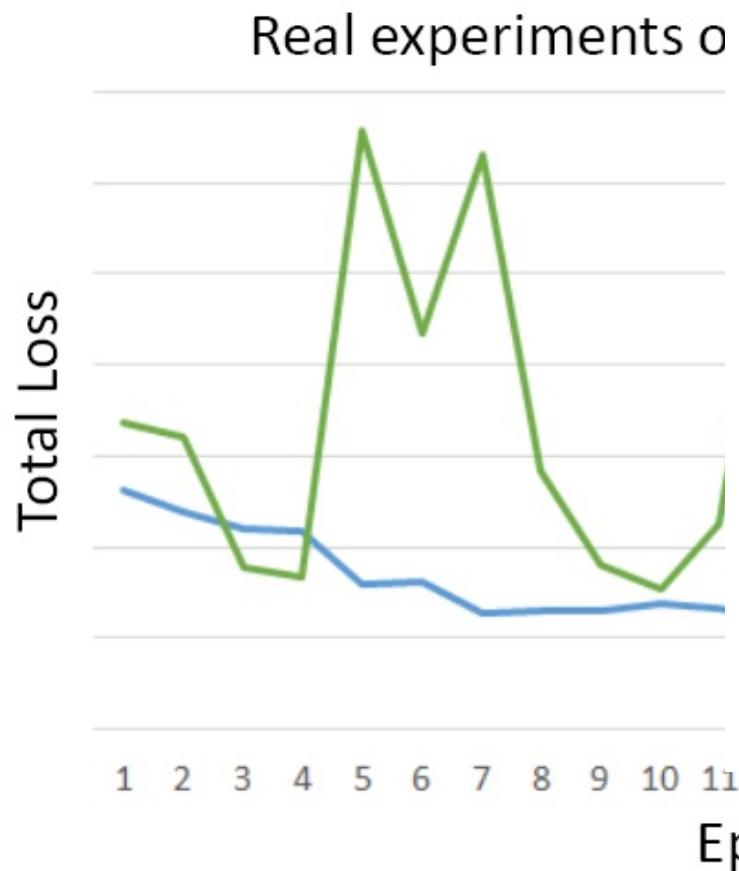


Learning



Unfortunately.....

- RNN-based network is not always easy to learn



Real experiments on RNN



 更多圖片

Tomáš Mikolov

電腦科學家

譯自英文 - Tomáš Mikolov is a Czech computer scientist working in the field of machine learning. He currently works at the Institute of Information, Robotics and Mechatronics of the German Aerospace Center (DLR). Mikolov has made significant contributions to deep learning and natural language processing. He is known for inventing the Word2vec word embedding method.

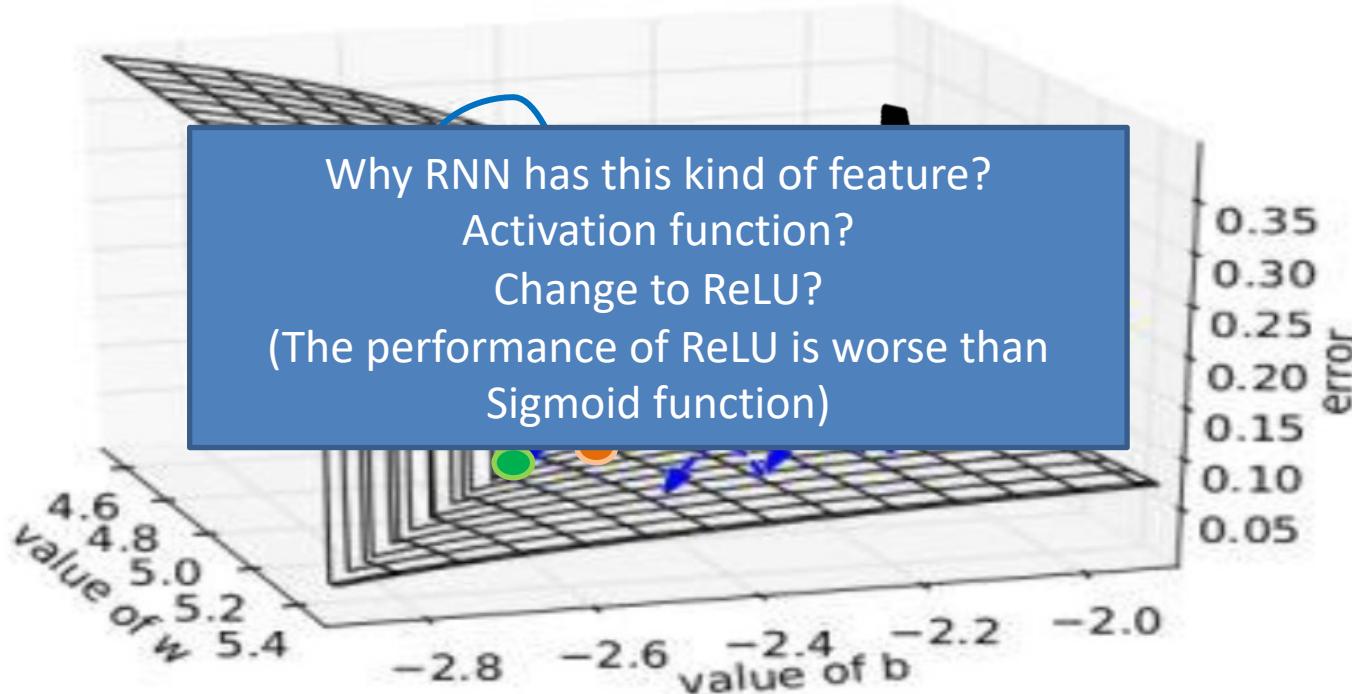
[維基百科 \(英文\)](#)

[查看原文說明 ▾](#)

The Error Surface is Rough

- Error surface: Total loss vs. Parameter changes
- Very flat or very steep on error surface in practice.

Clipping: When gradient is larger than threshold, gradient equals to current value and stop.



Why?

$$\begin{array}{ll} w = 1 & \rightarrow y^{1000} = 1 \\ w = 1.01 & \rightarrow y^{1000} \approx 20000 \end{array}$$

Large
 $\partial L / \partial w$

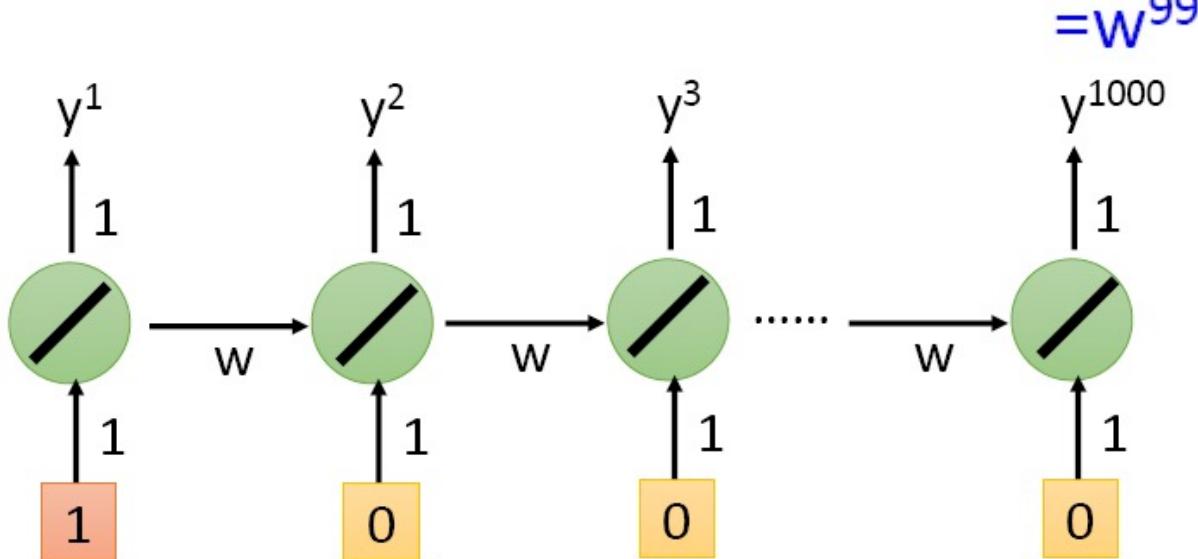
Small
Learning rate?

$$\begin{array}{ll} w = 0.99 & \rightarrow y^{1000} \approx 0 \\ w = 0.01 & \rightarrow y^{1000} \approx 0 \end{array}$$

small
 $\partial L / \partial w$

Large
Learning rate?

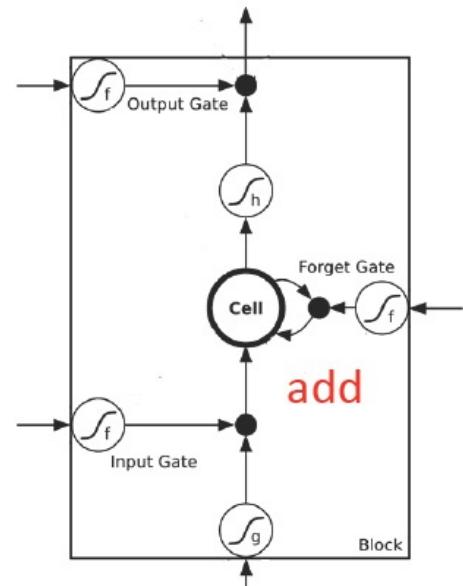
Toy Example



Helpful Techniques

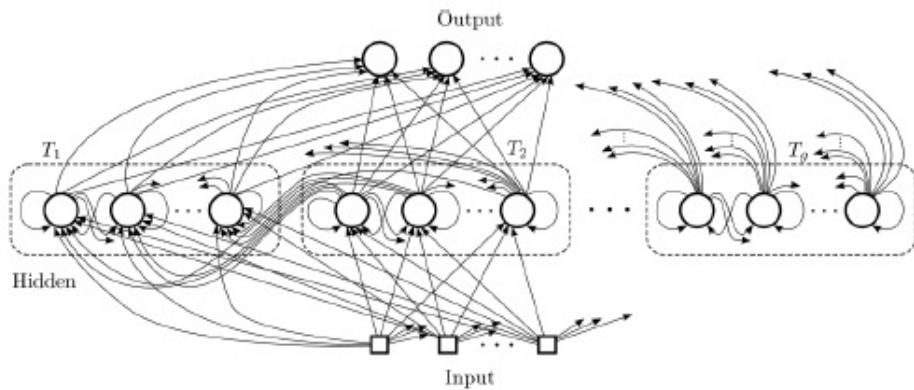
- Long Short-term Memory (LSTM)
 - Can deal with gradient vanishing (not gradient explode)
 - Memory and input are added
 - The influence never disappears unless forget gate is closed.
-  No Gradient vanishing
(If forget gate is opened.)

Gated Recurrent Unit(GRU): simpler than LSTM.
Training is more robust than LSTM.



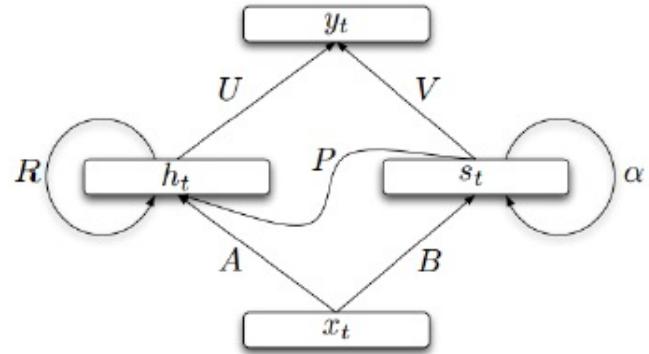
Helpful Techniques

Clockwise RNN



[Jan Koutnik, JMLR'14]

Structurally Constrained
Recurrent Network (SCRN)



[Tomas Mikolov, ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]

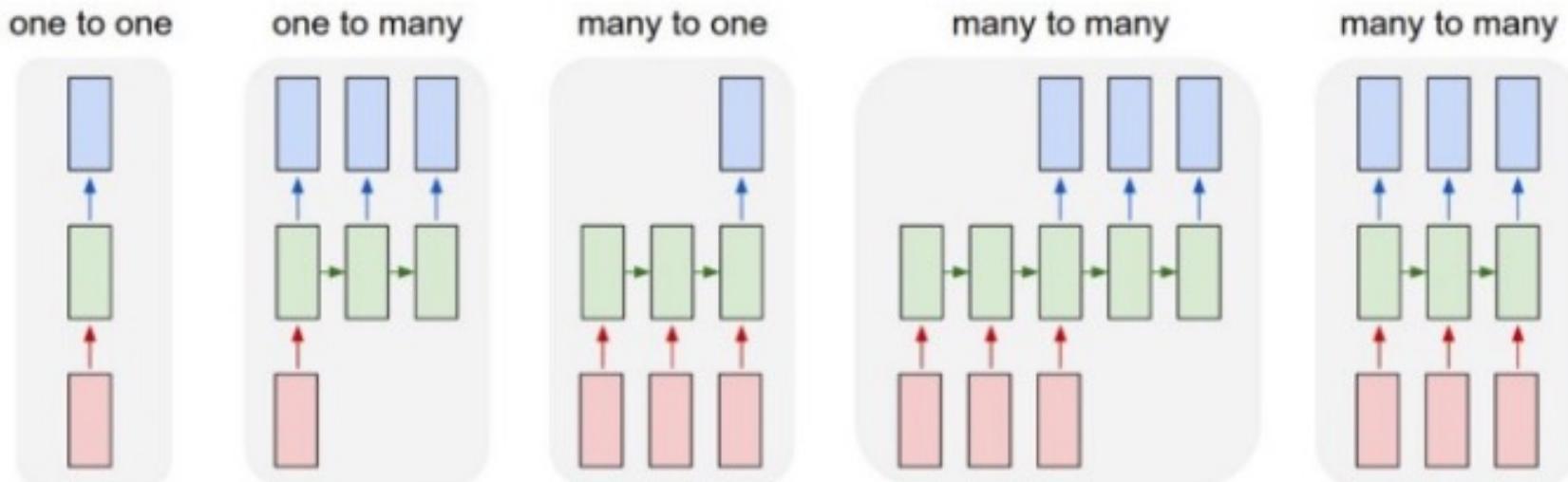
- Outperform or be comparable with LSTM in 4 different tasks



More Applications of RNN

Multiple Models

- The number of input and output of RNN can be different.

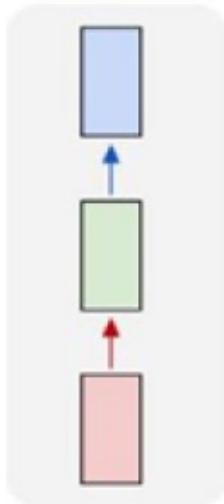


EX: POS Tagging

One to One

- Classification task
- DNN

Cat Dog





Many to One

- Input is a vector sequence, but output is only one vector

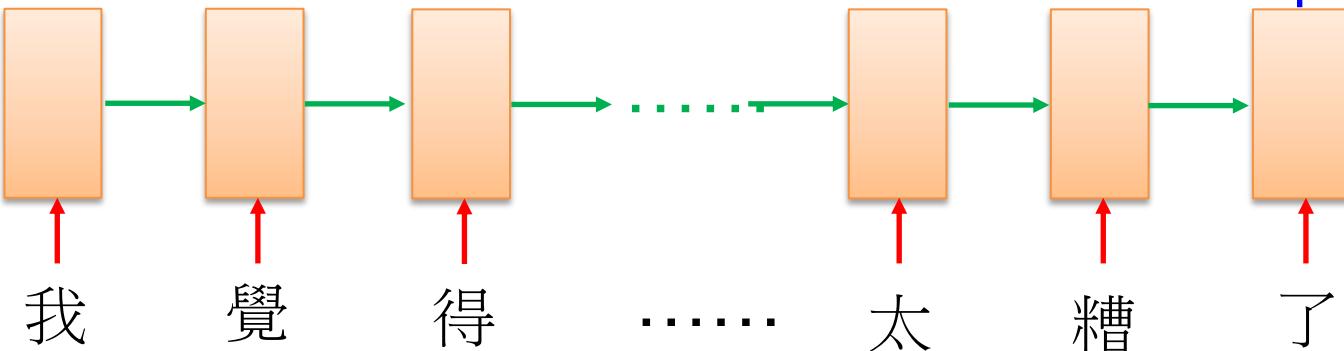
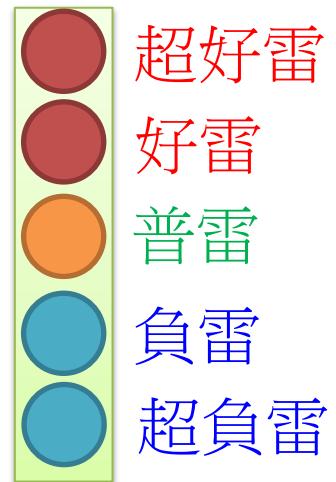
Sentiment Analysis

看了這部電影覺
得很高興

這部電影太糟了
.....

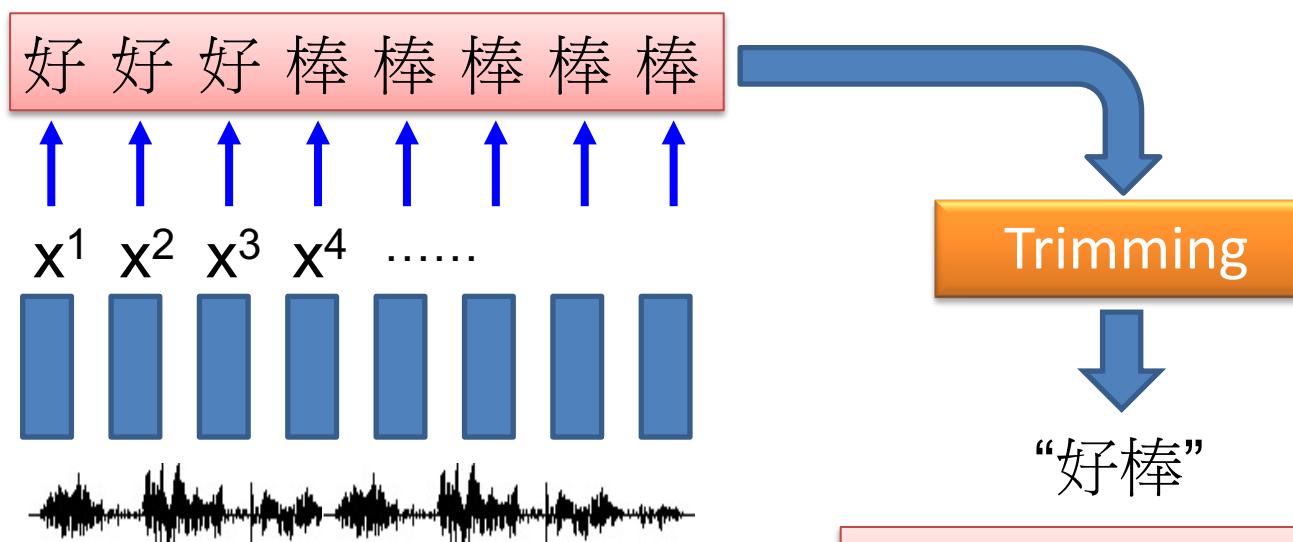
這部電影很棒

Positive (正雷) Negative (負雷) Positive (正雷)



Many to Many (Output is shorter)

- Both input and output are vector sequences, **but the output is shorter.**



Speech Recognition

You can never
recognize “好棒棒” !

Many to Many (Output is shorter)



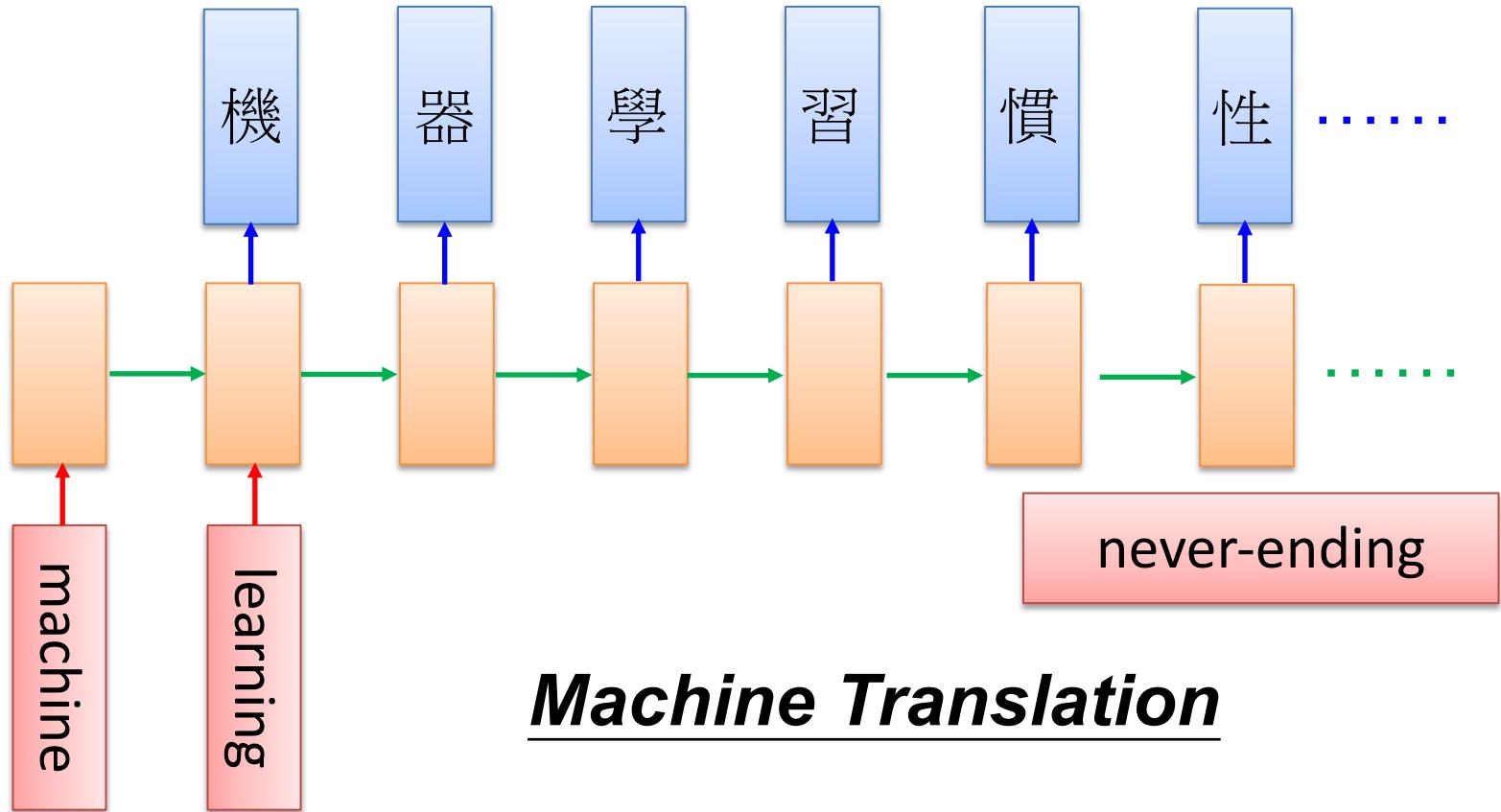
- Both input and output are vector sequences, **but the output is shorter.**
 - Connectionist Temporal Classification (CTC)
 - Add an extra symbol “ ϕ ” (同上)

好 φ φ 棒 φ φ φ φ → “好棒”

好 φ φ 棒 φ 棒 φ φ → “好棒棒”

Many to Many (No Limitation)

- Both input and output are vector sequences with different lengths. → Sequence to sequence learning





Many to Many (No Limitation)

- 推文接龍
 - Ref: <http://pttpedia.pixnet.net/blog/post/168133002-%E6%8E%A5%E9%BE%8D%E6%8E%A8%E6%96%87>

推xxx: ptt萬歲

推dd: 歲平安

噓ddf: 全

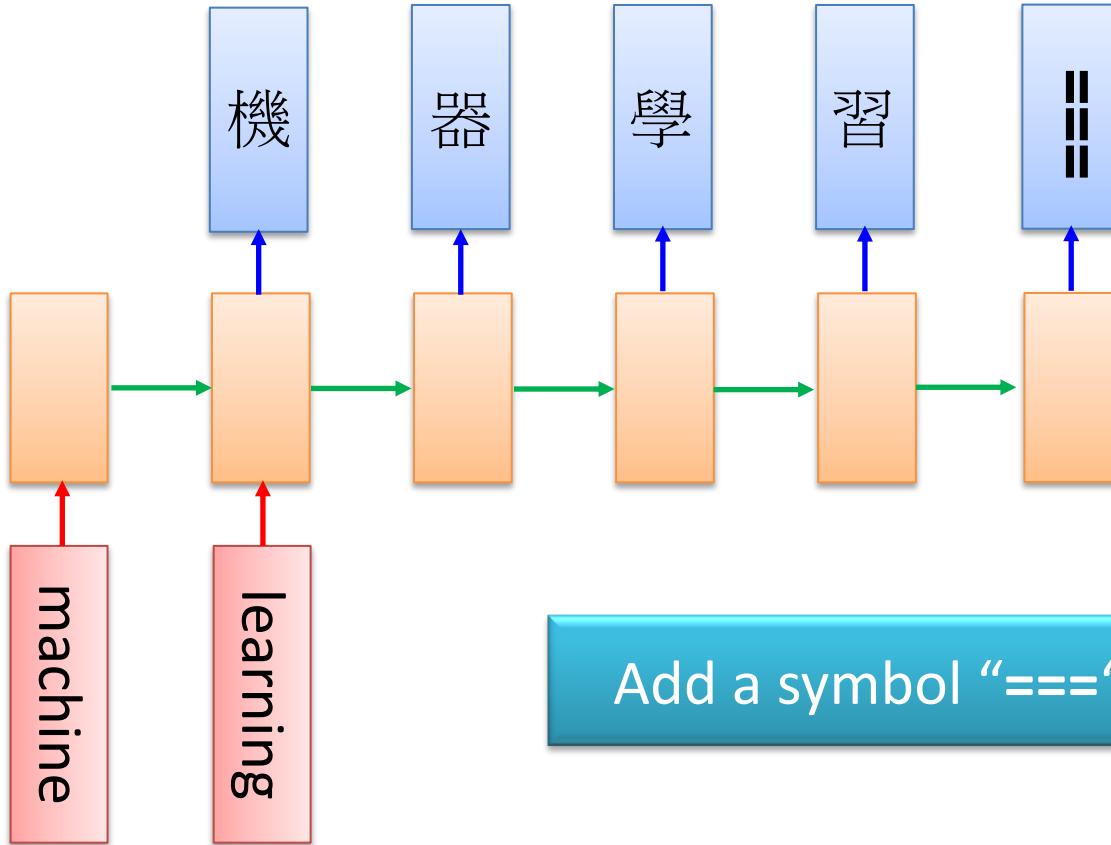
推zzzzzzzzzz: 家就是你家

: :

推tlkagk: =====斷=====

Many to Many (No Limitation)

- Both input and output are vector sequences with different lengths. → Sequence to sequence learning



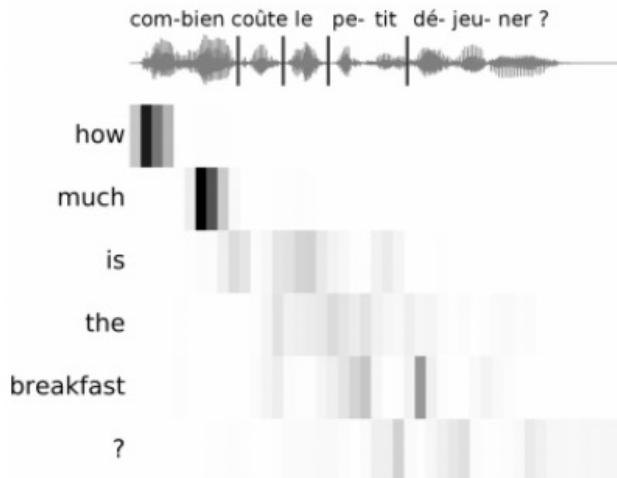
Many to Many (No Limitation)

- Both input and output are both sequences with different lengths. ->Sequence to sequence learning
- EX: Machine Translation (English translates to Chinese).

Is there possible input sounds of English and output the text of Chinese without speech recognition?



(a) Machine translation alignment

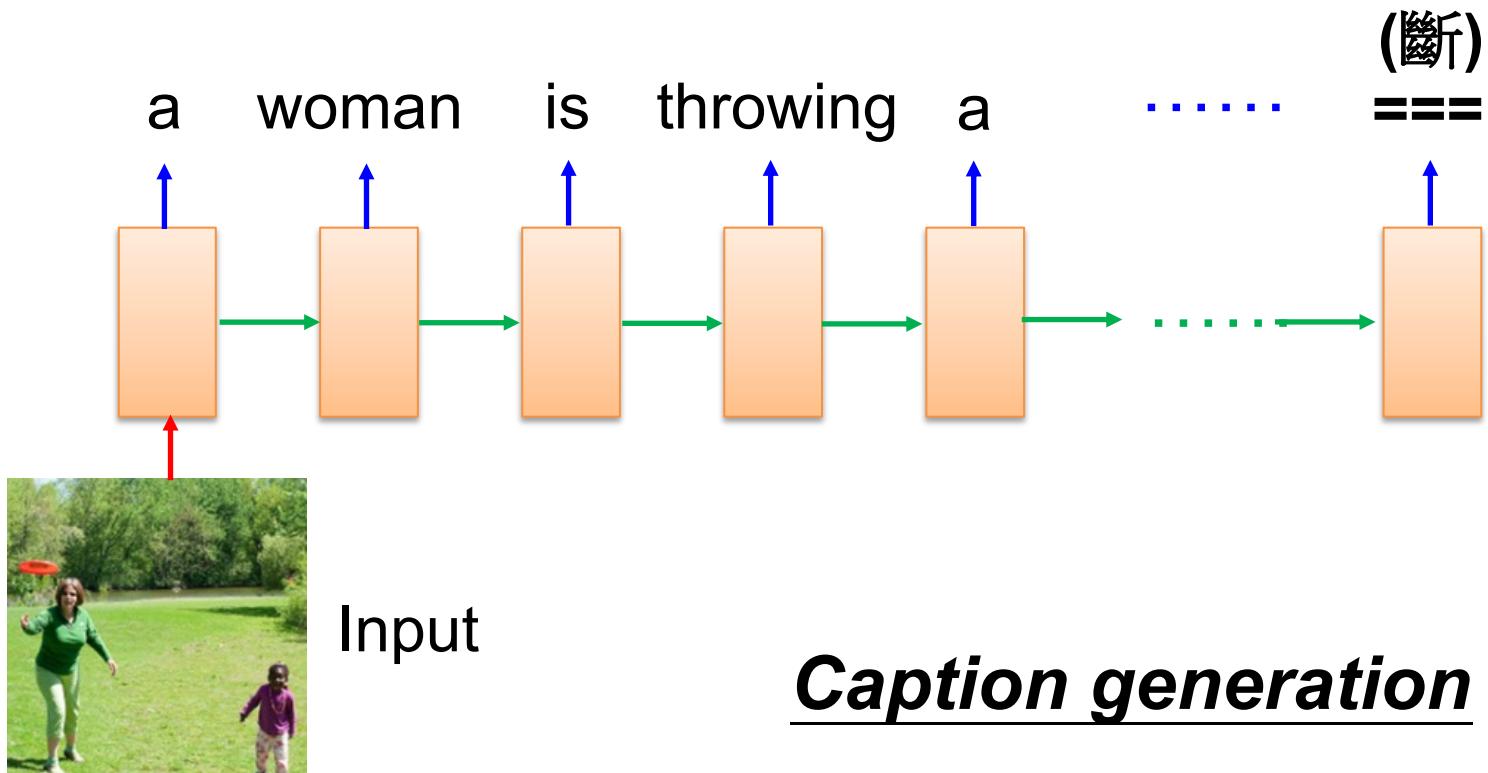


(b) Speech translation alignment

Figure 1: Alignments performed by the attention model during training

One to Many

- Input is one vector, but output is a vector sequence





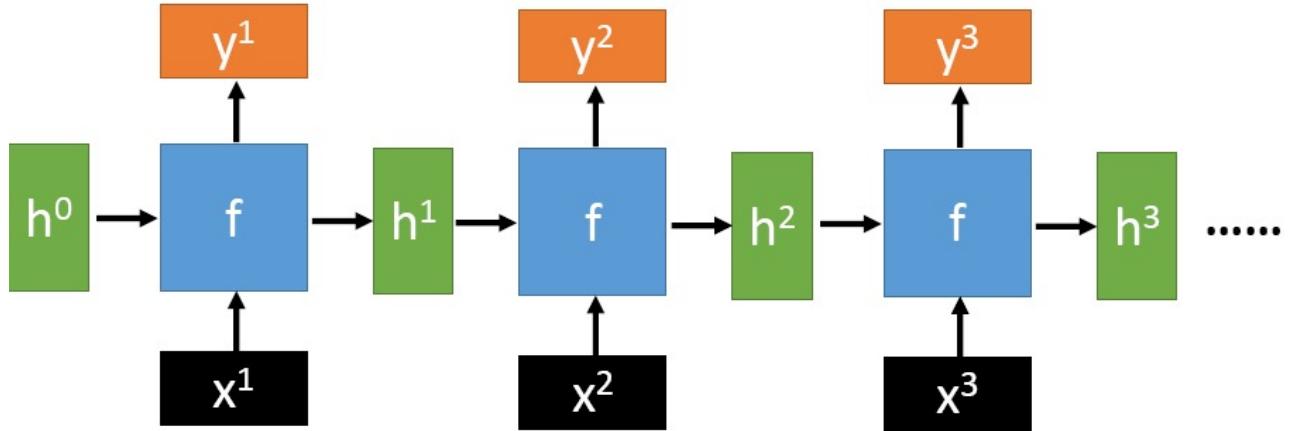
RNN for Anomaly Detection



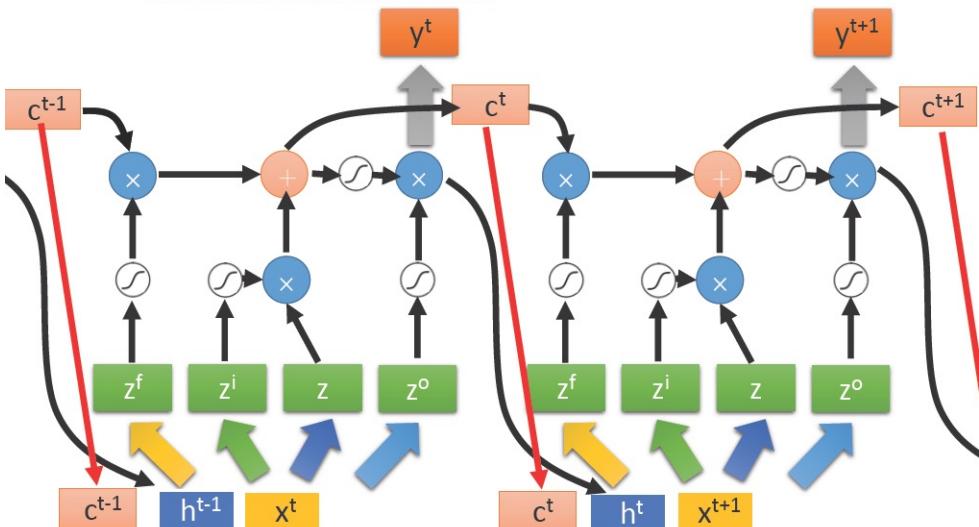
Summary

Summary

- Recurrent Neural Network



- LSTM





Recommended Reading List

- The Unreasonable Effectiveness of Recurrent Neural Networks
 - <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- Understanding LSTM Networks
 - <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>